

Scaling carbon fluxes from eddy covariance sites to globe: Synthesis and evaluation of the FLUXCOM approach

Martin Jung¹, Christopher Schwalm², Mirco Migliavacca¹, Sophia Walther¹, Gustau Camps-Valls³, Sujan Koirala¹, Peter Anthoni⁴, Simon Besnard^{1,5}, Paul Bodesheim^{1,6}, Nuno Carvalhais^{1,7}, Frédéric Chevallier⁸, Fabian Gans¹, Daniel S. Goll⁹, Vanessa Haverd¹⁰, Philipp Koehler¹¹, Kazuhito Ichii^{12,13}, Atul K. Jain¹⁴, Junzhi Liu^{1,15}, Danica Lombardozzi¹⁶, Julia E.M.S. Nabel¹⁷, Jacob A. Nelson¹, Michael O'Sullivan¹⁸, Martijn Pallandt¹⁹, Dario Papale^{20,21}, Wouter Peters²², Julia Pongratz^{23,17}, Christian Rödenbeck¹⁹, Stephen Sitch¹⁸, Gianluca Tramontana^{20,3}, Anthony Walker²⁴, Ulrich Weber¹, Markus Reichstein¹

¹Department of Biogeochemical Integration, Max Planck Institute for Biogeochemistry, Jena, 07745, Germany

²Woods Hole Research Center, Falmouth, MA, 02540-1644, USA

³Image Processing Laboratory (IPL), Universitat de València, Paterna, 46980, Spain

⁴Institute of Meteorology and Climate Research – Atmospheric Environmental Research (IMK-IFU), Karlsruhe Institute of Technology, Garmisch-Partenkirchen, 82467, Germany

⁵Laboratory of Geo-Information Science and Remote Sensing, Wageningen University and Research, Wageningen, 6708 PB, Netherlands

⁶Department of Mathematics and Computer Science, Friedrich-Schiller Universität Jena, Jena, 07743, Germany

⁷Departamento de Ciências e Engenharia do Ambiente (DCEA), Faculdade de Ciências e Tecnologia, FCT, Universidade Nova de Lisboa, Caparica, 2829-516, Portugal

⁸Laboratoire des Sciences du Climat et de l'Environnement (LSCE/IPSL), Université Paris-Saclay, Gif-sur-Yvette, F-91198, France

⁹Department of Geography, University of Augsburg, Augsburg, 86159, Germany

¹⁰Department Continental Biogeochemical Cycles, CSIRO Oceans and Atmosphere, Canberra, 2601, Australia

¹¹Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, USA

¹²Center for Environmental Remote Sensing (CEReS), Chiba University, Chiba, 263-8522, Japan

¹³Center for Global Environmental Research, National Institute for Environmental Studies, Tsukuba, 305-8506, Japan

¹⁴Department of Atmospheric Science, University of Illinois, Urbana, IL 61801, USA

¹⁵School of Geography, Nanjing Normal University, Nanjing, 210023, China

¹⁶Climate and Global Dynamics Laboratory, National Center for Atmospheric Research, Boulder, CO 80307, USA

¹⁷Department Land in the Earth System (LES), Max Planck Institute for Meteorology, Hamburg, 20146, Germany

¹⁸College of Life and Environmental Sciences, University of Exeter, Exeter, EX4 4QE, UK

¹⁹Department of Biogeochemical Systems, Max Planck Institute for Biogeochemistry, Jena, 07745, Germany

²⁰Department of Innovation in Biology, Agri-food and Forest systems (DIBAF), University of Tuscia, Viterbo, 01100, Italy

²¹Impacts on Agriculture, Forests and Ecosystem Services (IAFES), EuroMediterranean Center on Climate Change (CMCC), Lecce, 01100, Italy

²²Department of Meteorology and Air Quality, Wageningen University and Research, Wageningen, 6700 AA, Netherlands

²³Department of Geography, Ludwig-Maximilians-Universität München, München, 80333, Germany

²⁴Climate Change Science Institute, Oak Ridge National Laboratory, Oak Ridge, USA

Correspondence to: Martin Jung (mjung@bgc-jena.mpg.de)

51 **Abstract.** FLUXNET comprises globally-distributed eddy covariance-based estimates of carbon fluxes between
52 the biosphere and the atmosphere. Since eddy covariance flux towers have a relatively small footprint and are
53 distributed unevenly across the world, upscaling the observations is necessary to obtain global-scale estimates of
54 biosphere-atmosphere exchange. Based on cross-consistency checks with atmospheric inversions, sun-induced
55 fluorescence (SIF) and dynamic global vegetation models (DGVM), we provide here a systematic assessment of
56 the latest upscaling efforts for gross primary production (GPP) and net ecosystem exchange (NEE) of the
57 FLUXCOM initiative, where different machine learning methods, forcing datasets, and sets of predictor
58 variables were employed.

59 Spatial patterns of mean GPP are consistent across FLUXCOM and DGVM ensembles ($R^2 > 0.94$ at 1° spatial
60 resolution) while the majority of DGVMs show, for 70% of the land surface, values outside the FLUXCOM
61 range. Global mean GPP magnitudes for 2008-2010 from FLUXCOM members vary within 106 and 130 PgC yr⁻¹
62 with the largest uncertainty in the tropics. Seasonal variations of independent SIF estimates agree better with
63 FLUXCOM GPP (mean global pixel-wise $R^2 \sim 0.75$) than with GPP from DGVMs (mean global pixel-wise $R^2 \sim$
64 0.6). Seasonal variations of FLUXCOM NEE show good consistency with atmospheric inversion-based net land
65 carbon fluxes, particularly for temperate and boreal regions ($R^2 > 0.92$). Interannual variability of global NEE in
66 FLUXCOM is underestimated compared to inversions and DGVMs. The FLUXCOM version which uses also
67 meteorological inputs shows a strong co-variation of interannual patterns with inversions ($R^2 = 0.87$ for 2001-
68 2010). Mean regional NEE from FLUXCOM shows larger uptake than inversion and DGVM-based estimates,
69 particularly in the tropics with discrepancies of up to several hundred gC m² yr⁻¹. These discrepancies can only
70 partly be reconciled by carbon loss pathways that are implicit in inversions but not captured by the flux tower
71 measurements such as carbon emissions from fires and water bodies. We hypothesize that a combination of
72 systematic biases in the underlying eddy covariance data, in particular in tall tropical forests, and a lack of site-
73 history effects on NEE in FLUXCOM are likely responsible for the too strong tropical carbon sink estimated by
74 FLUXCOM. Furthermore, as FLUXCOM does not account for CO₂ fertilization effects carbon flux trends are
75 not realistic. Overall, current FLUXCOM estimates of mean annual and seasonal cycles of GPP as well as
76 seasonal NEE variations provide useful constraints of global carbon cycling, while interannual variability
77 patterns from FLUXCOM are valuable but require cautious interpretation. Exploring the diversity of Earth
78 Observation data and of machine learning concepts along with improved quality and quantity of flux tower
79 measurements will facilitate further improvements of the FLUXCOM approach overall.

80 **1 Introduction**

81 Upscaling local eddy covariance (EC) measurements (Baldocchi et al., 2001) from tower footprint to global
82 wall-to-wall maps uses globally-available predictor variables such as satellite remote sensing and meteorological
83 data (Jung et al., 2011). This forcing data is first used to establish empirical models for fluxes of interest at site
84 level, and then to estimate gridded fluxes by applying these models across all vegetated grid cells. Previous
85 FLUXNET upscaling efforts using machine learning techniques (Beer et al., 2010; Jung et al., 2009; Jung et al.,
86 2011) yielded global products that present a data-driven ‘bottom-up’ perspective on carbon fluxes between the
87 biosphere and the atmosphere. These ‘bottom-up’ products are complementary to process-based model
88 simulations and ‘top-down’ atmospheric inversions. However, estimates of carbon fluxes are subject to
89 uncertainty from choice of machine learning algorithm and predictor variables, forcing data, FLUXNET

90 measurements and incomplete representation of the different ecosystems therein. The FLUXCOM initiative
91 (www.fluxcom.org) aims to improve our understanding of the multiple sources and facets of uncertainties in
92 empirical upscaling and, ultimately, to provide an ensemble of machine learning-based global flux products to
93 the scientific community. Within FLUXCOM an intercomparison was conducted for two complementary
94 experimental setups of input drivers and resulting global gridded products. These setups systematically vary
95 machine learning and flux partitioning methods as well as forcing datasets to separate measured net ecosystem
96 exchange (NEE) into gross primary productivity (GPP) and Terrestrial Ecosystem Respiration (TER) (Jung et
97 al., 2019; Tramontana et al., 2016).

98
99 Evaluating the strengths and weaknesses of the FLUXCOM products and the approaches used therein is crucial
100 to inform potential scientific uses, and to guide future methodological developments. An evaluation based on
101 site-level cross-validation analysis (Tramontana et al., 2016) showed a general high consistency among machine
102 learning algorithms, experimental setups and flux partitioning methods applied in FLUXCOM. However, the
103 conclusions from site-level cross-validation may be limited by potential systematic measurement errors that are
104 inherent in the underlying EC measurements (e.g. Aubinet et al., 2012), or the spatially biased distribution of
105 FLUXNET sites (Papale et al., 2015). Therefore, cross-consistency checks of the FLUXCOM products with
106 independent estimates are important to consider. But such checks are complex due to limitations of the
107 independent approaches or the lack of comparability of similar but not identical variables. In this study, we
108 contextualize FLUXCOM products in relation to independent state-of-the-art estimates of carbon cycling. The
109 comparison strategy prioritises robust features of the independent datasets, and discusses residual uncertainties.

110
111 The objectives of this paper are (1) to present a synthesis and evaluation of FLUXCOM ensembles for GPP and
112 NEE against patterns of remotely sensed sun induced fluorescence (SIF) and atmospheric inversion results
113 respectively, (2) to discuss limitations of FLUXCOM and synthesize lessons learned, and (3) to outline potential
114 future paths of FLUXCOM development. Due to limitations of the SIF product with respect to interannual
115 variability (Zhang et al., 2018), the evaluation of GPP against SIF is restricted to seasonal variations of
116 photosynthesis. To reduce the impact of atmospheric transport-related uncertainties of inversion products, mean
117 annual and seasonal variations of NEE are compared at regional scales while interannual variability is assessed
118 at global scale. In addition, we contextualize our comparisons with FLUXCOM by providing comparisons with
119 the previous Model Tree Ensemble (MTE) results of Jung et al., 2011 (Ju11) as well as an ensemble of process-
120 based Global Dynamic Vegetation Model (DGVM) simulations from the TRENDY DGVM Projects (Le Quéré
121 et al., 2018; Sitch et al., 2015). Even though FLUXCOM also produced global products of TER, these are not
122 shown here due to a lack of an independent observational benchmark.

123 **2 Data and methods**

124 **2.1 FLUXCOM**

125 We used the cross-validated and trained machine learning techniques for the FLUXCOM carbon fluxes of
126 Tramontana et al. (2016) and generated large ensembles ($n = 120$) of global gridded flux products for two
127 different setups: remote sensing (RS) and remote sensing plus meteorological/climate forcing (RS+METEO)
128 setups (Fig. 1). In the RS setup, fluxes are estimated exclusively from Moderate Resolution Imaging

129 Spectroradiometer (MODIS) satellite data. In RS+METEO, fluxes are estimated from mean seasonal cycles of
130 satellite data and daily meteorological information (see Table S1). For the rationale of these setups, we refer the
131 interested reader to Tramontana et al., 2016 and Jung et al., 2019. For the RS setup, nine machine learning
132 methods were used to generate gridded products at an 8-daily temporal and 0.0833° spatial resolution for the
133 2001-2015 period. For the RS+METEO setup, three machine learning methods with five global climate forcing
134 data sets (Table 1) yielded products with daily temporal and 0.5° spatial resolution and time periods depending
135 on the meteorological data. The meteorological data included WATCH Forcing Data-ERA Interim (WFDEI;
136 Weedon et al., 2014), Global Soil Wetness Project 3 forcing data (GSWP3, Kim, 2017), CRU-JRA version 1.1
137 (Harris, 2019), ERA5 ((C3S), 2017), and a combination of observation-based radiation from CERES (Doelling
138 et al., 2013) and precipitation from GPCP (Huffman et al., 2001) (CERES-GPCP) resampled to 0.5°. The wide
139 range of data sources from reanalysis to station measurements to satellite observation is intentional and is meant
140 to bracket potential uncertainties in meteorological forcing.

141

142 For GPP and TER, we additionally considered uncertainty from flux partitioning methods by propagating two
143 different variants, one based on night-time NEE data (Reichstein et al., 2005) and one on daytime data (Lasslop
144 et al., 2010). Within the RS and RS+METEO setups, we followed a full factorial design of machine learning
145 methods (9 for RS, 3 for RS+METEO), flux partitioning variants (2 for GPP and TER), and climate forcing
146 input products (5, only for RS+METEO). Descriptions of machine learning methods, training, and validation
147 setup are available in Tramontana et al., 2016. The methodology of generating the global products is documented
148 in detail in the overview paper on global energy fluxes from FLUXCOM (Jung et al., 2019).

149

150 To allow for a better reuse of the large archive, we generated ensemble products of monthly values where
151 individual ensemble members were first aggregated to monthly means (Fig. 1). The ensemble products
152 encompass estimates of different machine learning estimates, flux partitioning variants for GPP and TER, and
153 different climate input data for RS+METEO. For the RS+METEO setup, this was also done separately for each
154 climate forcing data to allow modellers to compare their simulations with the FLUXCOM ensemble product
155 driven by the same forcing. The ensemble products (hereafter referred as FLUXCOM-RS and FLUXCOM-
156 RS+METEO) were generated as the median over ensemble members for each grid cell and month. The
157 FLUXCOM-RS products are based on 9 ensemble members for NEE and on 18 for GPP and TER. The
158 FLUXCOM-RS+METEO is based on 15 ensemble members for NEE and on 30 for GPP and TER.

159 **2.2 Process-model simulations (TRENDY)**

160 Dynamic Global Vegetation Models (DGVMs) represent an independent, process-based and bottom-up approach
161 to represent the terrestrial carbon cycle and its evolution with changing environmental conditions. Here we use
162 data from an ensemble of 16 DGVMs that were forced with the same climate (CRU-JRA v1.1), global
163 atmospheric CO₂ concentration, and land-use and land cover change data (S3 simulation) over the period 1700 –
164 2017, following a common protocol (TRENDY-v7) (Le Quéré et al., 2018; Sitch et al., 2015). This ensemble
165 provides fluxes at a monthly temporal resolution harmonized to a common 1° spatial resolution with simulations
166 from: CABLE-POP, CLASS-CTEM, CLM5.0, DLEM, ISAM, JSBACH, JULES, LPJ-GUESS, LPJ, OCN,
167 ORCHIDEE-CNP, ORCHIDEE-Trunk, SDGVM, SURFEX and VISIT. TER was calculated as the sum of
168 heterotrophic and autotrophic respiration; NEE as heterotrophic respiration minus net primary productivity. NBP

169 from models incorporates additional fluxes as well: fire emissions (10 DGVMs), land use change (all DGVMs),
170 harvest (14 DGVMs), grazing (6 DGVMs), and any other carbon flux in/out of the ecosystem (e.g. erosion, 1
171 DGVM, VISIT). LPJ-GUESS was excluded from comparisons of NEE or NBP since monthly output on
172 heterotrophic respiration was not available.

173 **2.3 Independent observation-based products**

174 For the comparison with GPP, we used gridded monthly SIF GOME-2 (Köhler et al., 2015) retrievals from the
175 far-red spectral range, and for the evaluation of NEE atmospheric inversion-based estimates from Jena
176 CarboScope (Rödenbeck et al., 2018), CAMSv17r1 (Chevallier et al., 2005; Chevallier et al., 2019), and
177 CarbonTracker-EU (CTE2018, Peters et al., 2010; van der Laan-Luijkx et al., 2017). We further include
178 comparisons to the previous GPP and NEE upscaling products of Jung et al., 2011 (hereafter referred as Ju11).

179 **2.4 Comparison approach**

180 **2.4.1 General considerations**

181 All products were harmonized to a common 1° spatial resolution with monthly temporal resolution as a basis of
182 all comparisons shown here. Cross-consistency checks for mean annual and mean seasonal variations of GPP
183 and NEE are based on the three year period 2008-2010. The time period is constrained by the availability of
184 GOME-2 data starting in 2008 and the corresponding end year of the RS+METEO ensemble with the GSWP3
185 forcing ending in 2010. The NEE interannual variability was initially assessed for 2001-2010 which is the
186 common period of the RS and RS+METEO ensembles while comparisons for longer-time periods were also
187 facilitated by using meteorological forcing specific RS+METEO products that cover longer time periods (Table
188 1).

189
190 FLUXCOM-RS and FLUXCOM-RS+METEO products are evaluated mostly separately. We report estimates for
191 the respective ensemble product (see section 2.1): the spread over individual ensemble members for uncertainty
192 and the mean of the ensemble members; the latter can be different from the ensemble product estimate (see
193 Sect.2.1). Occasionally, we use the range of estimates from the union of RS and RS+METEO ensemble
194 members to show the full FLUXCOM uncertainty range across the two setups (labelled as “FLUXCOM” only).
195 For the comparison of regional or global flux values, we used flux densities rather than integrated fluxes due to
196 inconsistencies in land-sea masks in different products. A common mask of valid data from the intersection of
197 FLUXCOM, TRENDY, and Ju11 was applied to all data streams, and a land area-weighted regional or global
198 mean calculated. Globally integrated GPP was calculated by scaling the global mean GPP density flux with the
199 global non-barren land area (122.4 Mio km²) derived from the MODIS land cover product (Friedl et al., 2010).
200 All reported R² values are squared Pearson’s correlation coefficients but negative correlation signs are
201 maintained through by multiplying R² values by -1. We aimed at structuring the cross-consistency checks with
202 SIF and inversion data to minimize confounding factors and uncertainties of the independent data that may have
203 affected the conclusions otherwise.

204 **2.4.2 Rationale of GPP-SIF comparison**

205 As the GPP-SIF relationship is approximately linear over seasonal time scales (Zhang et al., 2016), the
206 comparison was based on monthly values. To minimize confounding effects of canopy structure (e.g.
207 Migliavacca et al., 2017), the comparisons were done over time when canopy structure changes relative to GPP
208 changes are expected to be much weaker than spatial changes. The unstable orbit of the MetOp-A satellite that
209 carries one of the GOME-2 instruments and sensor degradation effects do not permit conclusive comparisons
210 with respect to interannual variability (Zhang et al., 2018). Therefore, we restricted the analysis to mean seasonal
211 cycles and show 1° maps of the R^2 between mean monthly GPP and SIF.

212

213 There are remaining caveats and uncertainties associated with the GPP-SIF relationship (see e.g. Porcar-Castell
214 et al., 2014 for an overview). Nevertheless, various studies have shown that SIF is currently the best proxy for
215 photosynthesis that can be remotely-sensed directly, in particular at seasonal time scale and over regions with
216 strong seasonal cycles. This is supported by strong empirical relationships between GPP and SIF across different
217 satellites and retrieval methods as well as from EC data, crop inventories, and data-driven GPP methods
218 (Frankenberg et al., 2011; Guanter et al., 2014; Joiner et al., 2018; Sun et al., 2017; Walther et al., 2016). This
219 gives us confidence in using SIF as an independent data stream for photosynthesis to evaluate FLUXCOM
220 products.

221 **2.4.3 Rationale of comparing net carbon fluxes with atmospheric inversions**

222 We compared atmospheric inversion-based net carbon release with FLUXCOM mean NEE at the seasonal scale
223 over the established 11 TRANSCOM regions (see Fig.S1 for a map) as atmospheric inversions are better
224 constrained over large spatial scales (Peylin et al., 2013). The comparison of interannual variability was
225 conducted at global scale due to its smaller signal and larger transport uncertainties compared to the seasonal
226 cycle. Due to various inversion uncertainties related to choices of atmospheric transport model, atmospheric
227 station CO₂ data, fossil fuel information, prior constraints, driving wind fields, and inversion strategy, we used
228 three different products: Jena CarboScope (s99oc_v4.3, Rödenbeck et al., 2018), CAMSv17r1 (Chevallier et al.,
229 2005; Chevallier et al., 2019), and CarbonTracker-EU (CTE2018, Peters et al., 2010; van der Laan-Luijkx et al.,
230 2017). To evaluate global NEE interannual variability patterns for periods since the late 1950s until present, we
231 further use two long-term atmospheric inversions (CarboScope s57Xoc_v4.3, sEXTocNEET_v4.3, Rödenbeck et
232 al., 2018) and annual CO₂ growth rate from the Global Carbon Budget (Le Quéré et al., 2018).

233

234 It is important to note that FLUXCOM NEE is semantically different from inversion-based net carbon exchange
235 between land and atmosphere. The former is solely the difference between gross fluxes (i.e., $NEE = TER - GPP$)
236 while the latter integrates all vertical movement of CO₂ including, for example, fire emissions, evasion from
237 inland waters, respired harvests, or volatile organic compounds (Kirschbaum et al., 2019; Zscheischler et al.,
238 2017). Simulations from TRENDY models report both, NEE and net biome productivity (NBP) which is
239 conceptually close but not identical to what atmospheric inversions provide. To assess whether conclusions are
240 affected by the different NEE vs NBP definitions we a) provide NEE and NBP estimates from TRENDY
241 models, b) we include comparisons where inversions were corrected for fire emissions (from CarbonTracker-
242 EU) to yield estimates closer to NEE, and c) discuss whether discrepancies with FLUXCOM can originate from
243 the omission of secondary carbon loss pathways given in the literature.

244 **3 Results and discussion**

245 **3.1 Gross primary productivity**

246 **3.1.1 Mean annual gross primary productivity**

247 Overall, our results suggest a high degree of cross-product (and, for FLUXCOM, also within-product)
248 consistency of global mean GPP patterns (Fig. 2). In fact, global patterns of mean GPP are consistent across both
249 FLUXCOM ensembles ($R^2=0.97$) as well as for Ju11 and TRENDY ensemble mean ($R^2>0.94$), despite sizeable
250 regional differences. The slope of the pair-wise 1:1 regressions among the different mean GPP data sets varies
251 within ~10%. FLUXCOM-RS shows about 10-20% lower GPP than FLUXCOM-RS+METEO in the highly
252 productive tropics and some subtropical regions. Both FLUXCOM setups estimate larger GPP than Ju11 and
253 TRENDY in some semi-arid regions and about 5-15% lower GPP in some extratropical areas. Despite a sizeable
254 total range of mean GPP from all 48 FLUXCOM members, the majority of TRENDY models (at least 9 out of
255 16) fall outside the FLUXCOM range for about 70% of the land surface (Fig. 3).

256
257 The mean global GPP of FLUXCOM-RS (111 PgC yr⁻¹) is about 10% lower than RS+METEO (120 PgC yr⁻¹,
258 Fig. 4), which is largely driven by differences in the tropics (Fig. 2). The cross-validation analysis indicated an
259 underestimation of FLUXCOM-RS GPP in the tropics (Tramontana et al., 2016), which was confirmed by a grid
260 cell-to-site data comparison for the FLUXNET 2015 data (which were not used for machine learning training
261 here) (Joiner et al., 2018). The reasons for the on-average lower GPP of RS compared to RS+METEO require
262 further investigation. It is unlikely that the smaller RS GPP values are because this setup is exclusively based on
263 remote sensing, as global latent heat from RS was larger than Ju11 (Jung et al., 2019). It seems to be rather
264 related to the specifically different predictor sets between RS and RS+METEO. This indicates that future
265 FLUXCOM efforts should expand the ensemble with respect to predictor set diversity to better account for this
266 source of uncertainty in upscaling. Focussing on FLUXCOM-RS+METEO, its ensemble spread (108-130 PgC
267 yr⁻¹) is much smaller than the TRENDY-based global GPPs (83-172 PgC yr⁻¹), and is primarily due to
268 differences among machine learning methods rather than meteorological forcing data (Fig.S2).

269
270 Our results imply that the present FLUXNET upscaling approach does not agree with larger GPP values of 150-
271 175 PgC yr⁻¹ derived from an isotope-based study (Welp et al., 2011). It is possible that the FLUXNET upscaling
272 approach underestimates GPP of highly managed and fertilized crops (Guanter et al., 2014) but their effects on
273 global GPP biases seem small (Joiner et al., 2018). At FLUXNET sites night-time CO₂ advection and storage
274 could cause underestimation of night-time CO₂ fluxes (Aubinet et al., 2012; McHugh et al., 2017; van Gorsel et
275 al., 2009) and thus underestimate GPP using the night-time NEE flux partitioning method. On the contrary, it has
276 been suggested that FLUXNET GPP estimated from the night-time partitioning method (Reichstein et al., 2005)
277 is overestimated as it ignores the effects of light inhibition of leaf respiration (Keenan et al., 2019; Wehr et al.,
278 2016) by on average 7% across FLUXNET sites (Keenan et al., 2019). But it should be noted that this value may
279 not be globally representative due to sizeable variations between ecosystems and with leaf area. Further, we only
280 find a small difference of mean global GPP of <2 PgC for day-time (Lasslop et al., 2010) and night-time
281 (Reichstein et al., 2005) NEE partitioning. This suggests that neither CO₂ advection nor the light inhibition of
282 leaf respiration appear to generate sizeable biases of global GPP in FLUXCOM—a tendency likely encouraged
283 by the relatively strict quality control on the EC fluxes data (Tramontana et al., 2016). Furthermore, a

284 comparison of EC-based GPP with biometric GPP estimates across 18 globally distributed sites showed good
285 agreement and no significant bias (Campioli et al., 2016). A recent study using Carbonyl Sulfide (COS)-based
286 partitioning for four contrasting European sites also showed good agreement with standard EC-based GPP where
287 systematic differences for mean GPP were $< 5\%$ (Spielmann et al., 2019). Therefore, we currently have no
288 strong indication that systematic biases of FLUXNET GPP propagate to global FLUXCOM GPP. Nevertheless,
289 we need to acknowledge that global GPP is largely driven by the productivity in the tropics where flux towers
290 are scarce and may be particularly uncertain due to challenging logistic and micrometeorological conditions (Fu
291 et al., 2018).

292
293 Various remote sensing-based light use efficiency approaches, calibrated with flux tower data, yielded global
294 GPP estimates of 109 (Zhao et al., 2005), 111 ± 21 (Yuan et al., 2010), 108-119 (Yu et al., 2018), 122 ± 25 (Jiang
295 and Ryu, 2016), 132 ± 22 (Chen et al., 2012), and 140 PgC yr^{-1} (Joiner et al., 2018). A simple calibration of only
296 near-infrared reflectance (NIRv) to EC data suggested a global GPP of 131-163 PgC yr^{-1} (Badgley et al., 2019).
297 Studies that assimilated atmospheric CO_2 concentration data into process model simulations yielded slightly
298 higher values of 148 (Anav et al., 2015) and $146 \pm 19 \text{ PgC yr}^{-1}$ (Koffi et al., 2012) with the latter study unable to
299 distinguish their best estimate from a global GPP of 117 PgC yr^{-1} because the atmospheric CO_2 alone cannot
300 constrain magnitudes of gross fluxes well. Assimilating SIF into process-models yielded 137 ± 6 (Norton et al.,
301 2019) and $166 \pm 10 \text{ PgC yr}^{-1}$ (MacBean et al., 2018). More recent isotope studies derived global GPP as 120 ± 30
302 PgC yr^{-1} (Liang et al., 2017), and global NPP of $\sim 60 \text{ PgC yr}^{-1}$ (Hellevang and Aagaard, 2015) which implies
303 global GPP of 109-150 PgC yr^{-1} considering a range of NPP:GPP ratios of 0.4-0.55. In conclusion, global
304 FLUXCOM GPP estimates are within the currently most plausible 110-150 PgC yr^{-1} range.

305 **3.1.2 Seasonal cycles of gross primary productivity**

306 Cross-consistency analysis of mean monthly GPP seasonal cycles from FLUXCOM with SIF from GOME-2
307 (Köhler et al., 2015) shows widespread and strong agreement for both FLUXCOM setups (Fig. 5), except for the
308 inner tropics where seasonality is weak and SIF retrievals might be affected by the South Atlantic Magnetic
309 Anomaly (Köhler et al., 2015). FLUXCOM-RS tends to show better agreement with SIF than FLUXCOM-
310 RS+METEO in agricultural regions of Southeast Asia, maybe because only the mean seasonal cycles of
311 remotely sensed land surface properties were used in the latter. Conversely, FLUXCOM RS+METEO shows on
312 average better consistency with SIF in some semi-arid regions, e.g., Australia. However, maps of the maximum
313 R^2 with SIF for RS and RS+METEO respectively have similar patterns with good agreement of both products in
314 Australia, and even in the tropics (Fig.S4). This suggests that the inclusion of some machine learning methods
315 somewhat negatively impacts the ensemble, especially for RS which shows larger spread (see Fig.S4 for mean
316 R^2 of the RS ensemble members). With SIF, both FLUXCOM setups show similar consistency as Ju11. The
317 consistency of FLUXCOM with SIF is much better than with TRENDY models, in particular in tropical and
318 subtropical regions. This implies that, despite sporadic spatial coverage of FLUXNET sites and previously
319 identified incomplete capturing of water stress (Bodesheim et al., 2018; Tramontana et al., 2016), FLUXCOM
320 still has a large potential to inform and constrain process-based model simulations of seasonal variations of
321 photosynthesis in moisture-limited regions.

322 **3.2 Net ecosystem exchange**

323 **3.2.1 Mean annual net ecosystem exchange**

324 In most TRANSCOM regions, FLUXCOM shows a stronger mean annual net carbon uptake than indicated by
325 atmospheric inversions with a particularly large systematic difference in the tropics (Fig. 6). This pattern of a
326 large tropical carbon sink in FLUXCOM is qualitatively consistent among the different FLUXCOM setups and
327 ensemble members, as well as with previous estimates from Ju11. To date, this is a systematic feature of the
328 current data-driven approach of upscaling EC measurements with machine learning.

329
330 Multiple independent approaches indeed imply a sizeable carbon sink in intact tropical forests (Arneth et al.,
331 2017; Gaubert et al., 2019; Pan et al., 2011), which appears to be largely or entirely offset by carbon loss
332 pathways in the tropical region such as fire, land-use change emissions, and evasion from inland waters. These
333 CO₂ sources are not sampled by EC measurements from FLUXNET, and are, therefore, not represented in
334 FLUXCOM. However, the missing fluxes only resolve up to roughly half of the gap (Zscheischler et al., 2017).
335 The comparatively small differences between net carbon release estimates by inversions and those where fire
336 emissions were corrected for, as well as the small differences between NEE and -NBP from TRENDY further
337 suggest that these secondary carbon loss fluxes do not drive the large discrepancy between FLUXCOM and
338 inversion-based mean net carbon exchange. Nevertheless, substantial uncertainty remains in the magnitude of
339 these secondary carbon fluxes and their incomplete accounting in TRENDY models and inversions (Kirschbaum
340 et al., 2019; Zscheischler et al., 2017).

341
342 Issues with the current FLUXCOM approach certainly contribute, likely dominate, the discrepancy between
343 atmospheric top-down and FLUXCOM mean NEE. Potential factors that could contribute to this are (1) a
344 FLUXNET sampling bias (see also Sect. 4.1.2) towards ecosystems with a large carbon sink, particularly in the
345 tropics (Saleska et al., 2003); combined with (2) missing predictor variables related to disturbance and site-
346 history (Amiro et al., 2010; Besnard et al., 2018, see also Sect. 4.2.1), or (3) biases of eddy covariance NEE
347 measurements, e.g. due to night-time advection of CO₂ (Hayek et al., 2018; van Gorsel et al., 2008), especially
348 under tall tropical forest canopies (Hutyra et al., 2008, Fu et al., 2018). Fu et al. (2018) studied 63 site-years of
349 EC data from 13 tropical forest sites and report a mean between-site NEE of -567 gC m⁻² yr⁻¹ showing that the
350 large tropical sink in FLUXCOM is inherited from FLUXNET data. The authors pointed out that for about half
351 of the sites where measurements of CO₂ concentration along the vertical profile were available and the storage
352 was considered in the NEE processing, the carbon sink was less than half (-340 gC m⁻² yr⁻¹) compared to those
353 without storage correction (-832 gC m⁻² yr⁻¹). However, the small sample size together with the large between-
354 site standard deviation of mean NEE (459 gC m⁻² yr⁻¹) not only makes robust conclusions difficult, but also
355 indicates potentially large diversity between tropical ecosystems. Clearly, more tropical EC sites are needed
356 along with a better accounting of systematic errors in EC-based NEE measurements to resolve this issue.

357 **3.2.2 Seasonal cycles of net ecosystem exchange**

358 We find a good consistency between FLUXCOM and inversions with respect to amplitude and shape of the
359 seasonal cycles of NEE in many TRANSCOM regions, especially over the North American Boreal, North
360 American Temperate, and Europe regions with R² values > 0.92 (Fig. 7). As with mean annual NEE, the

361 seasonal cycle mismatch relative to inversions may be linked to carbon loss fluxes not accounted for in
362 FLUXCOM, such as fire emissions that are seasonally relevant in tropical and subtropical regions. However,
363 adjusting inversion-based NBP towards NEE by correcting for fire emissions does not improve the
364 correspondence with FLUXCOM in tropical and subtropical regions (Fig.S5). In tropical regions, the weak
365 seasonality paired with comparatively large spread among inversions does not allow for robust conclusions.
366 Overall, the seasonal variations of FLUXCOM NEE show potential to constrain the large uncertainty in
367 TRENDY models, and potentially even atmospheric inversions at the regional scale, especially considering that
368 their uncertainty range across only three products is still significant.

369 **3.2.3 Interannual variability of net ecosystem exchange**

370 Spatial patterns of the magnitude of the interannual variability (IAV) of land carbon sink for the period 2001-
371 2010 share some common features among atmospheric inversions, FLUXCOM-RS, FLUXCOM-RS+METEO
372 and TRENDY. For example, all products identify the hotspots in southeast Asia, southern North America, and
373 also in the Siberian tundra (Fig. 8). Overall, there are still differences in the spatial patterns of IAV magnitude
374 among and within different data-streams.

375
376 All EC data-driven methods, in particular FLUXCOM-RS+METEO, underestimate magnitude of IAV compared
377 to inversions (Fig. 8). The reasons for the underestimation of IAV magnitude by FLUXCOM are not fully clear.
378 Within FLUXCOM, the smaller IAV magnitude of RS+METEO NEE compared to that of RS is linked to the
379 use of only mean seasonal cycles of RS-based land surface properties in RS+METEO setup. The IAV of carbon
380 loss fluxes that are not captured by FLUXCOM, such as through fire, are currently thought to be comparatively
381 small at the global scale and appear minor here (see Fig.S6). Machine learning methods already underestimate
382 the IAV at the site level (Marcolla et al., 2017; Tramontana et al., 2016). The low bias in FLUXCOM IAV is a
383 direct consequence of the comparatively small explained variance for NEE anomalies. Thus, improving the
384 predictability of NEE IAV at site level has potential to also correct the magnitude of globally integrated IAV
385 variance.

386
387 Despite the tendency of FLUXCOM products to underestimate IAV magnitude, FLUXCOM-RS+METEO
388 reproduces year-to-year variations of globally integrated annual land carbon exchange anomalies derived from
389 atmospheric inversions for 2001-2010 ($R^2=0.87$). It shows better consistency than TRENDY with one of the
390 long-term inversions (Fig.S7). Further examination of this ensemble reveals that the choice of machine learning
391 method, rather than meteorological forcing data, has a larger influence on IAV of global NEE (Fig.S8). Here, the
392 Random Forests method performed less well compared to the other two methods. Interestingly, training Random
393 Forests with an almost identical predictor set but at half-hourly temporal scale rather than at daily scale
394 (Bodesheim et al., 2018) substantially improved the R^2 (from 0.31 to 0.60, S8). This indicates that machine
395 learning methods can benefit from higher temporal variability provided by millions of high-frequency NEE
396 measurements, especially for signals such as IAV that are small and difficult to extract. In addition, underlying
397 functional relationships can be better extracted from high-frequency data as the predictor space is better covered,
398 allowing for improved discrimination of drivers that have stronger covariation on longer time-scales.

399

400 To better understand the qualitatively different global NEE IAV patterns between RS and RS+METEO setups,
401 we infer which NEE IAV signals are consistent or lacking among FLUXCOM setups and TRENDY models by
402 assessing correlation patterns (Fig. 9). We find the strongest consistencies of NEE IAV between FLUXCOM-RS
403 and FLUXCOM-RS+METEO in many semi-arid regions, and almost no consistency otherwise. This suggests
404 that the main discrepancies of globally integrated NEE IAV between FLUXCOM-RS and FLUXCOM-
405 RS+METEO are likely not due to differences in their capabilities of reflecting water stress effects. It has been
406 shown that despite the local dominance, water-related NEE anomalies largely cancel spatially in RS+METEO
407 and TRENDY resulting in the dominance of temperature-related NEE anomalies in globally integrated land sink
408 IAV (Jung et al., 2017, but see Humphrey et al., 2018 for a different perspective). Studies on effects of water
409 availability on spatial GPP anomalies using the RS data yielded highly plausible patterns that were consistent
410 with independent data (Flach et al., 2018; Orth et al., 2019; Walther et al., 2019). Also the comparison of
411 FLUXCOM-RS GPP monthly anomalies with the independent FLUXNET2015 data set showed unexpected
412 large consistency when anomalies were scaled by the site-specific observational range (Joiner et al., 2018).
413 When delineating the regions with larger agreement between RS+METEO and TRENDY than that between RS
414 and TRENDY, we can infer that FLUXCOM-RS seems to miss important NEE anomaly features in the tropics.
415 This is likely due to (1) a combination of sparse satellite data availability, cloud contamination, and geometrical
416 illumination effects in the tropics or (2) that the processes governing NEE IAV in the tropics cannot be captured
417 by satellite-based predictors alone in RS (even under ideal observational conditions) but require additional
418 meteorological variables such as temperature that is included in the RS+METEO setup. Some support for the
419 latter point comes from Byrne et al., 2019 who found strong correlations of anomalies from GOSAT inversions
420 with NEE from RS+METEO and soil temperatures in the tropics but not with SIF and a drought indicator,
421 suggesting that temperature impacts respiration more than photosynthesis in the tropics.

422
423 Overall there are large discrepancies among FLUXCOM and TRENDY as well as amongst TRENDY models
424 with respect to local NEE IAV. This reflects our limited understanding and capabilities to model year-to-year
425 variations of local ecosystem carbon exchange. Both data-driven and process-based approaches also showed
426 poor performance with respect to NEE IAV in FLUXNET sites (Tramontana et al., 2016, Morales et al., 2005).
427 However, both approaches yield good correspondence of globally integrated NEE with atmospherically-derived
428 interannual land sink variations. This correspondence is due to two reasons: first, the spatial compensation of
429 locally important processes that are not well captured by the models; and second, models capture better the
430 temperature-related signals that gain relevance at larger spatial scales (Jung et al., 2017). Whether the large
431 uncertainty of modelling NEE IAV at ecosystem level is due to misspecified parameterizations, missing
432 predictors, inaccurate forcing data and/or absent processes remains a research priority. Our understanding and
433 ability to model NEE IAV bottom-up would greatly benefit from atmospheric inversions that could localize NEE
434 robustly. Exploiting the massive space-based column CO₂ data in the future will hopefully facilitate the
435 improvements on this aspect. Despite large uncertainties and apparent knowledge gaps in NEE IAV from both an
436 observational and modelling perspective, there are promising indications of improved capability to track IAV
437 patterns with FLUXCOM such as the good correspondence of RS+METEO with inversions at global scale, and
438 independent verifications of GPP IAV of RS at least outside the wet tropics (Flach et al., 2018; Joiner et al.,
439 2018; Orth et al., 2019; Walther et al., 2019).

440 **4 Methodological limitations and potential ways forward**

441 Machine learning methods can learn arbitrarily complex functions and provide a nearly perfect model of a
442 phenomenon if they are fed with the right data and trained thoroughly. Thus the quality, quantity, and
443 completeness of the input data determine the quality of the output. In the following, we discuss the relevance of
444 limitations associated with data from the FLUXNET network, and of the limited capabilities of representing all
445 relevant factors by observable predictor variables. We also outline potential strategies for improvements, both
446 overall and with respect to machine learning approaches specifically. The continued and rapid development of
447 machine learning notwithstanding, we believe that the FLUXCOM approach is at present more limited by
448 available “information” rather than by available machine learning methods.

449 **4.1 FLUXNET observations**

450 **4.1.1 Potential observation errors**

451 The comparatively large random errors of high-frequency EC measurements diminish quickly when aggregated
452 to daily or 8-daily averages used here. Furthermore, training on half-hourly EC data (Bodesheim et al., 2018)
453 helps machine learning methods extract patterns from noisy data. In general, poor signal-to-noise ratios can be
454 counteracted by larger sample size. More problematic than random errors are potential systematic errors of EC
455 measurements since those would propagate to the derived global carbon flux products. Even though there have
456 been large efforts by the community to characterize and to correct for systematic errors, such as those due to low
457 turbulence and CO₂ advection (e.g. Aubinet et al., 2005; Aubinet et al., 2012; Papale et al., 2006), uncertainties
458 remain on the relevance and magnitude of those errors in the processed FLUXNET data. Differences due to
459 instrumentation and maintenance pose another potential source of uncertainty. Additionally, the energy balance
460 closure gap at FLUXNET sites is still not resolved (Stoy et al., 2013), while it remains unclear to what extent
461 this is relevant for CO₂ fluxes (Leuning et al., 2012). Systematic errors in GPP and TER derived from the flux
462 partitioning method of NEE based on night-time data (Reichstein et al., 2005) may arise due to the neglected
463 effect of inhibited photorespiration during daytime (Keenan et al., 2019; Wehr et al., 2016). Nevertheless, all
464 these issues together seem to be relatively small compared to the predominant patterns of variability in EC data,
465 e.g., seasonal variations, that are very consistent across FLUXCOM and independent observation-based data
466 streams shown here. The relatively strict quality controls on the flux training data (Tramontana et al., 2016) may
467 have been instrumental here. The trade-off between data quality and training data volume was not explicitly
468 studied in FLUXCOM, and related experimental setups would be desirable to gauge the robustness of the global
469 products shown here. Even small systematic errors in EC data could degrade important signals such as
470 interannual variability, trends, annual sums of NEE, or subtle differences between sites related to functional
471 properties (e.g., radiation use efficiency). Systematic errors that would be prevalent across the network would
472 result in systematic biases of derived global fluxes. For global GPP and energy fluxes (Jung et al., 2019), the
473 values obtained from FLUXCOM are generally consistent with current knowledge but our ability to
474 independently quantify such fluxes is also limited.

475 **4.1.2 Potential representation issues**

476 Ideally, a measurement network samples all relevant gradients of the driving factors and magnitudes of the
477 predicted quantities. There are several potential issues with the current sampling by FLUXNET sites. With

478 respect to relevance for net carbon exchange, there are carbon loss pathways that FLUXNET does not capture
479 such as fire emissions, CO₂ evasion from inland waters, and lateral exports due to harvest or erosion that are
480 respired elsewhere (Kirschbaum et al., 2019). The effects of strongly enhanced respiration in the years after large
481 disturbances (Amiro et al., 2010) are challenging to capture due to stochastic and destructive nature of
482 disturbances.

483

484 To meet the assumptions of EC method, FLUXNET stations are confined to reasonably flat terrain. Topographic
485 effects on ecosystem fluxes are primarily due to their influence on environmental drivers, i.e., the predictor
486 variables. Thus, the extrapolation to hillslopes should be reasonable if the topographic effects are accounted for
487 in the gridded predictor variables. This might be challenging especially for remote sensing products due to
488 necessary but complicated corrections of illumination conditions. The uncertainties of these topographic factors
489 might become particularly relevant and should be studied for prediction of fluxes at a higher spatial resolution.
490 For the current FLUXCOM products with rather coarse spatial resolution, we expect that topographic effects are
491 reflected in the predictor variables and the remaining subpixel heterogeneity largely cancel out.

492

493 Perhaps the most fundamental and frequent critique of the FLUXNET upscaling approach is related to the
494 spatially clumped geographic distribution of EC sites in North America, Europe, Japan, and now Australia with
495 only sparsely distributed towers elsewhere (Schimel et al., 2015). However, what matters eventually for machine
496 learning methods is how well the predictor space, rather than geographic space, is sampled. To assess this, we
497 developed an extrapolation index (EI) that estimates the expected additional relative error of a flux prediction
498 due to a large distance to the nearest training data in the predictor space (S2). We applied this method for GPP
499 and FLUXCOM-RS training data as an example, and found that the conditions that are least well represented by
500 FLUXNET are associated to primarily extremely cold and dry regions (Figure 10). Surprisingly, the humid
501 tropics are well represented in the predictor space suggesting that the environmental conditions represented by
502 the predictor set are well sampled by the data from FLUXNET sites. The extremely cold and dry conditions that
503 seem to constitute the biggest extrapolation issues are typically associated with small GPP fluxes and thus also
504 small prediction errors. To account for that, we spatialized the expected GPP error of the RS ensemble (Figure
505 10, see S2 for details), which largely scales with GPP magnitude but also shows patterns of larger expected
506 errors in semi-arid regions than that expected from flux magnitude alone. The multiplication of the expected
507 GPP error with the extrapolation index provides the extrapolation severity index (ESI) that shows where poor
508 FLUXNET sampling likely increases the absolute prediction error strongly. According to these results, sub-
509 tropical semi-arid regions, in particular India, appear as most affected, suggesting that GPP upscaling from
510 FLUXNET would benefit most strongly from improved data availability for towers representing these
511 conditions. Despite these limitations of data, we found excellent consistency of FLUXCOM GPP seasonal cycles
512 with SIF over these regions, which was in fact much better than the consistency between TRENDY models and
513 SIF. This suggests that while more towers in semi-arid regions will help reduce uncertainty in future upscaling
514 efforts, FLUXCOM can already provide useful information for constraining the models in these regions. It also
515 shows that the bias in geographic representation of FLUXNET sites is not as critical as anticipated due to the
516 flexibility and adaptiveness of machine learning methods. The sampled environmental conditions (predictor
517 space) should cover the conditions of the global application domain rather than being representative of it. The

518 larger issue of the FLUXNET representation bias is associated with drawing conclusions from the site-level
519 cross-validation because the evaluation metrics are easily biased towards certain regions and ecosystems.

520

521 The methodology used here to assess the extrapolation problem quantitatively has several limitations. For
522 example, potential differences in EC data quality were not accounted for. Perhaps, the largest but unavoidable
523 limitation is the reliance on the predictor set and the assumption that it captures all relevant gradients. In a sense,
524 the methodology can only uncover “known unknowns”. If an important predictor is missing, the method would,
525 of course, not see any extrapolation penalty with respect to the missing factor. Somewhat ironically, we may
526 need more towers in the first place to identify further relevant predictors in an objective way to, say, better
527 capture the diversity in the tropics (Fu et al., 2018) or in agricultural systems (Guanter et al., 2014) where we
528 anticipate that the current sampling is limiting the FLUXCOM approach.

529 **4.2 Driving factors and predictors**

530 Assuming infinite sample size, perfect quality and coverage, the success of machine learning methods depends
531 entirely on the completeness of the predictor set for the target variable, given an adequate training. The predictor
532 set for FLUXNET upscaling is practically constrained by 1) the availability of consistent observations at site
533 level across all sites, and for most of their temporal coverage at a spatial resolution sufficiently close to the flux
534 tower footprint; and 2) the availability of corresponding global grids at an adequate spatial and temporal
535 resolution and temporal coverage. This explains the predictor space of remotely sensed land products from
536 MODIS along with tower-measured meteorology chosen in FLUXCOM. While the general success of the
537 FLUXCOM approach suggests that the predictor sets contain sufficient information for predicting the variability
538 of carbon fluxes, it is also obvious that some factors are not well accounted for.

539 **4.2.1 Site-history**

540 It has been argued previously (Besnard et al., 2018; Jung et al., 2011; Tramontana et al., 2016) that the current
541 limitations of unrealistic mean NEE patterns from FLUXNET upscaling is also due to missing predictor
542 variables that describe site history effects such as forest age or time since disturbance. These factors have been
543 shown to influence IAV (Musavi et al., 2017; Tamrakar et al., 2018) and to drive mean NEE patterns in synthesis
544 studies (e.g. Amiro et al., 2010). Including forest age in a simple empirical model helped predicting between site
545 variations of mean NEE across FLUXNET sites (Besnard et al., 2018). Counterintuitively, including forest age
546 in training a machine learning method on monthly NEE did not improve the predictability of mean site NEE
547 (Besnard et al., 2019), albeit possibly due to data or methodological limitations. We find the largest
548 discrepancies of mean FLUXCOM NEE with atmospheric inversions in the tropics, where site history plays a
549 substantial role in NEE magnitude (Pugh et al., 2019), but the concept of forest age is hardly applicable due to
550 the generally uneven aged nature of stands, and reliable estimates of gridded age, e.g., from forest inventories are
551 not available. Efforts to incorporate the information from long-term LANDSAT time series to capture site
552 history effects did not reveal an improvement in the predictions of mean NEE, but it remains unclear if this was
553 due to limited information content in these time series or due to methodological issues (Besnard et al., 2019).
554 Thus, this issue remains a significant scientific challenge. Potentially, the availability and application of high-
555 resolution biomass and vegetation optical depth estimates from radar remote sensing along with a carefully

556 collected ancillary data on biomass, basal area, tree diameter and tree age distributions at ICOS and NEON sites
557 may help in the future.

558 **4.2.2 Management**

559 We are presumably lacking important information on anthropogenic management effects, in particular for crops
560 (Guanter et al., 2014) but also for forests. This is primarily due to a lack of information on, e.g., crop type,
561 fertilizer application, irrigation, harvest or thinning at FLUXNET sites, but also due to the still-limited number of
562 crop sites to provide sufficient information on relevant predictors therein. Accounting for the management
563 effects in the FLUXCOM approach either by explicit management information or implicitly by adequate remote
564 sensing data may also help improve the predictions of IAV of local-scale carbon fluxes, in particular with cross-
565 validation since most FLUXNET sites are subject to some degree of management.

566 **4.2.3 CO₂ fertilisation**

567 FLUXCOM lacks any explicit treatment of the effects of CO₂ fertilization causing carbon flux trends to be
568 unrealistic (Fig.S11). This is a challenging problem due to a comparatively small size of [CO₂] effect. This, in
569 turn, makes it particularly vulnerable to distortions through measurement uncertainties, and, on an annual scale,
570 largely indistinguishable from any other factor that varies with time. Potentially, in the future, the availability of
571 longer time series along with high-quality near surface atmospheric CO₂ data at high spatial and temporal
572 resolution at the tower scale could allow for extracting a CO₂ fertilization effect by exploiting diurnal, seasonal,
573 and spatial CO₂ gradients in addition to the long-term trend.

574 **4.2.4 Water stress**

575 Site-level cross-validation analysis (Bodesheim et al., 2018; Tramontana et al., 2016) indicated that soil moisture
576 effects on carbon fluxes are not always well captured. In RS+METEO, moisture effects are explicitly addressed
577 by a simple meteorology driven water availability index. The RS setup relies entirely on indirect information
578 encoded in remotely sensed surface properties such as vegetation indices and land surface temperatures. The
579 comparison of FLUXCOM GPP seasonal cycles with SIF yielded excellent agreement, also in water limited
580 systems, and studies on drought effects using the GPP RS product (Flach et al., 2018; Orth et al., 2019; Walther
581 et al., 2019) found plausible patterns that were consistent with independent data on large scales. Nevertheless,
582 we should strive further to improve water stress effects in the upscaling approach given its significance. Better or
583 explicit predictor variables on soil moisture may help. Unfortunately, current soil moisture products from remote
584 sensing are only representative of the top few centimeters and are at comparatively coarse spatial resolution
585 limiting their applicability in reflecting spatial heterogeneities of soil moisture. Perhaps, the larger issue is
586 diverse ecosystem specific responses to soil moisture variations due to different ecosystem compositions, rooting
587 patterns, plant hydraulics, stomata and other physiological traits. Thus, exploring remotely sensed products that
588 reflect additional or complementary information on water stress effects, such as diurnal cycles of land surface
589 temperature from geostationary satellites, is a potential way forward.

590 **4.2.4 Product properties**

591 The success of incorporating novel informative data of site properties in the FLUXCOM approach is always
592 contingent on the quality of the corresponding global gridded products. Systematic differences between a

593 predictor variable used for training at the site-level and global forcing data, as well as any potential artefacts due
594 to retrieval issues or merging different data records spatially or temporally propagate to global flux products.
595 Future improvements of the FLUXCOM approach will thus require progress in other research fields with
596 emphasis on the processing, correction, and harmonization of Earth observation products. Especially for
597 remotely sensed data, strategies to bridge scales of satellite pixels, overpass times, and repeat cycles to
598 continuous measurements of flux footprints are needed. In addition, making use of novel data in the FLUXCOM
599 framework requires the concurrent development of new methodological strategies to cope with the small
600 temporal overlap of the FLUXNET data history. More generally, the quality and quantity of Earth observation
601 data has been increasing rapidly, bringing challenges and opportunities for upscaling.

602 **4.3 Machine learning**

603 **4.3.1 Exploiting temporal data structures**

604 The machine learning methods employed in FLUXCOM are classic ones, while novel approaches could bring
605 further improvements. One conceptual limitation of all machine learning methods used in FLUXCOM is that
606 they assume independent and identically distributed (i.i.d.) variables, and thus do not respect or exploit temporal
607 structures in the training data. This problem can be remedied by using other machine learning methods based on
608 convolutions. For example, recurrent neural networks (RNNs) were designed for time-series and can account for
609 dynamics such as ecosystem lag and memory effects on carbon flux variability. Conceptually, lag and memory
610 effects emerge due to the effect of unobserved ecosystem state variables. RNNs can potentially counteract the
611 lack of a relevant state variable in the predictor set if the state variable's instantaneous effect is encoded in the
612 temporal history of other predictor variables (e.g., current soil moisture as a function of previous weather). While
613 exploiting the temporal information of predictors using an RNN improved predictions of monthly carbon fluxes
614 in terms of the seasonal cycle and thereby also across-site variability, predictions of interannual variability were
615 not improved as compared to exploiting only time-instantaneous effects based on site-level cross-validation
616 (Besnard et al., 2019). Further exploration of the machine learning methods that exploit the temporal structure of
617 predictors has a potential to improve FLUXCOM upscaling.

618 **4.3.2 Promising strategies**

619 Deep learning techniques, in general, and convolutional neural networks (CNNs), in particular, have proven to
620 be very powerful especially for image processing and recognition tasks (LeCun et al., 2015). Their conceptual
621 strength lies in the automated extraction of features, in particular those related to spatial structures that render the
622 design and implementation of hand-crafted predictor variables unnecessary. Whether simply employing CNNs
623 for upscaling brings similar improvements over traditional machine learning techniques as in other domains is
624 questionable. This is because the number and spatial distribution of FLUXNET towers seems insufficient to
625 exploit the power of CNNs to extract relevant features of spatial structure. However, combining CNNs with
626 transfer learning approaches seems very promising from a conceptual perspective. The principle of transfer
627 learning is to learn relevant features from a more densely observed proxy variable of the actual target and use the
628 feature representation for learning the target (Pan and Yang, 2010). The learning of the proxy variable can be
629 done either prior to or simultaneously with the actual target such that information from much larger sample of
630 the proxy can be transferred to the sparsely observed target variable. This approach could be applicable to the
631 upscaling of FLUXNET GPP by using remotely sensed SIF as a proxy and thereby alleviate issues related to

632 small sample size (e.g., extrapolation) but also aid the extraction of small but relevant signals (e.g., IAV). Spatial
633 structures in high-resolution SIF data may further encode effects of management or topographically controlled
634 soil moisture variations that could be exploited with CNNs and improve predictions.

635
636 Hybrid approaches, i.e. the integration of machine learning method with process understanding and physical
637 constraints, are another promising avenue. This allows for different strategies and levels of complexity are
638 possible (Reichstein et al., 2019), and could also greatly help in regularizing machine learning predictions to be
639 sensible under extrapolation conditions. In the context of FLUXCOM, for, say, constraining the anticipated weak
640 signal of CO₂ fertilization in observations within theoretically derived bounds, would allow this relevant yet
641 observationally poorly constrained dynamic to be incorporated. If the hybrid approach features the
642 conceptualization of fluxes and pools as in process models, it would also allow for constraints by multiple
643 complementary data streams simultaneously.

644
645 An important aspect to improve in the future is also the quantification of uncertainty in the predictions, including
646 the propagation of observational uncertainties. Gaussian processes are now computationally tractable for big data
647 problems while providing probabilistic confidence intervals and allowing for uncertainty propagation (Camps-
648 Valls et al., 2016; Wang et al., 2019). Combining Gaussian Processes with deep neural nets (You et al., 2017) or
649 designing deep Gaussian process models (Damianou and Lawrence, 2013) are powerful new machine learning
650 tools with the potential to improve FLUXCOM.

651 **Conclusions**

652 The FLUXCOM initiative generated a large ensemble of global carbon flux products for two defined setups that
653 differ in the set of predictor variables and spatial-temporal resolution. The ensemble is comprised of 120
654 products using up to 9 machine learning algorithms, two flux-partitioning variants for GPP and TER, and 5
655 meteorological forcing data sets. The large and systematically generated ensemble allows for assessing and
656 studying uncertainties of the fluxes as well as the approaches used in FLUXCOM. We assessed FLUXCOM
657 GPP and NEE patterns against remotely sensed sun-induced fluorescence (SIF), atmospheric inversions and
658 process model simulations from the TRENDY initiative.

659
660 We found strong consistency of FLUXCOM with SIF and atmospheric inversions with respect to seasonal
661 variations, highlighting FLUXCOM's suitability to evaluate and constrain seasonal cycles for processed-based
662 and top-down approaches. The global GPP from RS+METEO was 120 ± 7 PgC yr⁻¹ (mean \pm 1 s.d.), while the
663 global GPP from RS (111 ± 3 PgC yr⁻¹) is lower likely due to underestimation in the tropics. FLUXCOM shows a
664 consistently large carbon sink in the tropics that can, at present, not be reconciled with our knowledge derived
665 from atmospheric CO₂ constraints; possibly implying an underestimation of carbon loss and/or missing carbon
666 loss pathways by FLUXNET observations. Patterns of year-to-year variations of the global land carbon sink
667 from FLUXCOM-RS+METEO show good consistency with atmospheric inversions, while magnitudes of
668 interannual variability are underestimated in the data-driven approaches. As FLUXCOM lacks the effect of CO₂
669 fertilization, trends are not realistic and should only be used for assessing the exclusive effects of climate
670 changes on carbon fluxes.

671
672 Moving forward, increasing the size of the FLUXNET network, improving its quality, standardization and
673 coverage will both improve quality and reduce uncertainties in the upscaling approach. This holds especially
674 with respect to signals that are important but relatively small and difficult to extract such as interannual
675 variability or trends. Increasing the number of tropical sites alone would also help constrain global flux
676 magnitudes, and, in particular, would help resolve the large tropical carbon sink shown by FLUXCOM but
677 missing in atmospheric inversions. Based on the number of registered FLUXNET sites alone, an approximate
678 five-fold increase in the number of sites with available data seems feasible in theory; if all respective researchers
679 would contribute their flux data to the global community effort. This indicates that any efforts to improve eddy
680 covariance data, sharing, harmonization and processing are crucial.

681
682 Beyond extending the data frame, the current FLUXCOM intercomparison suggests that the next phase of
683 methodological developments should be to move away from predetermined setups and instead towards a set of
684 dedicated experiments that explore novel strategies of data integration with machine learning method (e.g., deep,
685 transfer, and hybrid approaches) and, more importantly, the diversity in the potential predictor space from Earth
686 Observation data. Within FLUXCOM, we find the largest differences between RS and RS+METEO setups
687 which primarily differ in the set of input predictor variables. Thus, the current approach of upscaling FLUXNET
688 measurements seems more information rather than algorithm limited.

689
690 Overall, the success of FLUXCOM approach depends on the interplay of many different factors. Monitoring our
691 progress will be essential but challenging, and must combine site-level cross-validation, cross-consistency
692 checks with global independent data-streams, novel and dedicated experiments as well as tailored validations of
693 methods with artificial data similar to Observation System Simulation Experiments. Despite the many
694 challenges, integrating eddy covariance ecosystem scale fluxes, Earth Observation data and machine learning
695 method has already proven valuable in many respects despite being a comparatively new field. An exciting and
696 challenging future lays ahead; that the contribution of experts in different fields combined with open and real
697 time data sharing could lead to a unique semi-operational carbon monitoring system. This in turn provides a
698 promising perspective to unify and synergistically exploit data-driven biospheric bottom-up and atmospheric
699 top-down approaches.

700 **Data availability**

701 Monthly carbon flux data of all ensemble members as well as the ensemble estimates from the FLUXCOM
702 initiative (<http://www.fluxcom.org>) are freely available (CC4.0 BY licence) from the data portal of Max Planck
703 Institute for Biogeochemistry (<https://www.bgc-jena.mpg.de/geodb/projects/Home.php>) after registration.
704 Choose 'FluxCom' in the dropdown menu of the database and select FileID 260. The users will be provided with
705 an access to an ftp server. The ftp directory is structured in a consistent way and stores files with consistent
706 naming convention in netcdf-4 format (see S3 for details).. Products with daily or 8-daily temporal resolution or
707 customized ensemble estimates are available on request to Martin Jung (mjung@bgc-jena.mpg.de). TRENDY
708 model output is available on request to Stephen Sitch (S.A.Sitch@exeter.ac.uk).

709 **Author contributions**

710 MJ conceived the study, performed the analysis, and drafted the manuscript with intellectual input and extensive
711 edits from all co-authors.

712 **Competing interests**

713 The authors declare no competing interests.

714 **Acknowledgements**

715 The authors acknowledge funding from European Space Agency Climate Change Initiative ESA-CCI RECCAP2
716 project (ESRIN/4000123002/18/I-NB), and EU H2020 projects, CHE (GA 776186), VERIFY (GA 776810), E-
717 SHAPE (GA 820852), and BACI (GA 640176). We further want to thank Ana Bastos for input on an earlier
718 version of the manuscript.

719 **References**

720 (C3S), C. C. C. S.: ERA5: Fifth generation of ECMWF atmospheric reanalyses of the global climate. (CDS), C.
721 C. C. S. C. D. S. (Ed.), 2017.

722 Amiro, B. D., Barr, A. G., Barr, J. G., Black, T. A., Bracho, R., Brown, M., Chen, J., Clark, K. L., Davis, K. J.,
723 Desai, A. R., Dore, S., Engel, V., Fuentes, J. D., Goldstein, A. H., Goulden, M. L., Kolb, T. E., Lavigne, M. B.,
724 Law, B. E., Margolis, H. A., Martin, T., McCaughey, J. H., Misson, L., Montes-Helu, M., Noormets, A.,
725 Randerson, J. T., Starr, G., and Xiao, J.: Ecosystem carbon dioxide fluxes after disturbance in forests of North
726 America, *Journal of Geophysical Research: Biogeosciences*, 115, 2010.

727 Anav, A., Friedlingstein, P., Beer, C., Ciais, P., Harper, A., Jones, C., Murray-Tortarolo, G., Papale, D., Parazoo,
728 N. C., Peylin, P., Piao, S., Sitch, S., Viovy, N., Wiltshire, A., and Zhao, M.: Spatiotemporal patterns of terrestrial
729 gross primary production: A review, *Reviews of Geophysics*, 53, 785-818, 2015.

730 Arneeth, A., Sitch, S., Pongratz, J., Stocker, B. D., Ciais, P., Poulter, B., Bayer, A. D., Bondeau, A., Calle, L.,
731 Chini, L. P., Gasser, T., Fader, M., Friedlingstein, P., Kato, E., Li, W., Lindeskog, M., Nabel, J. E. M. S., Pugh,
732 T. A. M., Robertson, E., Viovy, N., Yue, C., and Zaehle, S.: Historical carbon dioxide emissions caused by land-
733 use changes are possibly larger than assumed, *Nature Geoscience*, 10, 79, 2017.

734 Aubinet, M., Berbigier, P., Bernhofer, C. H., Cescatti, A., Feigenwinter, C., Granier, A., Grunwald, T. H.,
735 Havrankova, K., Heinesch, B., Longdoz, B., Marcolla, B., Montagnani, L., and Sedlak, P.: Comparing CO₂
736 storage and advection conditions at night at different carboeuroflux sites, *Boundary-Layer Meteorology*, 116, 63-
737 94, 2005.

738 Aubinet, M., Feigenwinter, C., Heinesch, B., Laffineur, Q., Papale, D., Reichstein, M., Rinne, J., and van Gorsel,
739 E.: Nighttime flux correction. In: *Eddy Covariance: A Practical Guide to Measurement and Data Analysis*,
740 Aubinet, M., Vesala, T., and Papale, D. (Eds.), Springer Atmospheric Sciences, Springer, Dordrecht, 2012.

741 Badgley, G., Anderegg, L. D. L., Berry, J. A., and Field, C. B.: Terrestrial gross primary production: Using
742 NIRV to scale from site to globe, *Global Change Biology*, 0, 2019.

743 Baldocchi, D., Falge, E., Gu, L. H., Olson, R., Hollinger, D., Running, S., Anthony, P., Bernhofer, C., Davis, K.,
744 Evans, R., Fuentes, J., Goldstein, A., Katul, G., Law, B., Lee, X. H., Malhi, Y., Meyers, T., Munger, W., Oechel,
745 W., U, K. T. P., Pilegaard, K., Schmid, H. P., Valentini, R., Verma, S., Vesala, T., Wilson, K., and Wofsy, S.:
746 FLUXNET: A new tool to study the temporal and spatial variability of ecosystem-scale carbon dioxide, water
747 vapor, and energy flux densities, *Bulletin of the American Meteorological Society*, 82, 2415-2434, 2001.

748 Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rödenbeck, C., Arain, M. A.,
749 Baldocchi, D., Bonan, G. B., Bondeau, A., Cescatti, A., Lasslop, G., Lindroth, A., Lomas, M., Luysaert, S.,
750 Margolis, H., Oleson, K. W., Rouspard, O., Veenendaal, E., Viovy, N., Williams, C., Woodward, F. I., and
751 Papale, D.: Terrestrial gross carbon dioxide uptake: global distribution and covariation with climate, *Science*,
752 329, 834-838, 2010.

753 Besnard, S., Carvalhais, N., Arain, M. A., Black, A., Brede, B., Buchmann, N., Chen, J., Clevers, J. G. P. W.,
754 Dutrieux, L. P., Gans, F., Herold, M., Jung, M., Kosugi, Y., Knohl, A., Law, B. E., Paul-Limoges, E., Lohila, A.,

755 Merbold, L., Rouspard, O., Valentini, R., Wolf, S., Zhang, X., and Reichstein, M.: Memory effects of climate
756 and vegetation affecting net ecosystem CO₂ fluxes in global forests, *PLOS ONE*, 14, e0211510, 2019.

757 Besnard, S., Carvalhais, N., Arain, M. A., Black, A., de Bruin, S., Buchmann, N., Cescatti, A., Chen, J., Clevers,
758 J. G. P. W., Desai, A. R., Gough, C. M., Havrankova, K., Herold, M., Hörtnagl, L., Jung, M., Knohl, A., Kruijt,
759 B., Krupkova, L., Law, B. E., Lindroth, A., Noormets, A., Rouspard, O., Steinbrecher, R., Varlagin, A., Vincke,
760 C., and Reichstein, M.: Quantifying the effect of forest age in annual net forest carbon balance, *Environmental
761 Research Letters*, 13, 124018, 2018.

762 Bodesheim, P., Jung, M., Gans, F., Mahecha, M. D., and Reichstein, M.: Upscaled diurnal cycles of land–
763 atmosphere fluxes: a new global half-hourly data product, *Earth Syst. Sci. Data*, 10, 1327-1365, 2018.

764 Byrne, B., Jones, D. B. A., Strong, K., Polavarapu, S. M., Harper, A. B., Baker, D. F., and Maksyutov, S.: On
765 what scales can GOSAT flux inversions constrain anomalies in terrestrial ecosystems?, *Atmos. Chem. Phys.
766 Discuss.*, 2019, 1-42, 2019.

767 Campioli, M., Malhi, Y., Vicca, S., Luysaert, S., Papale, D., Peñuelas, J., Reichstein, M., Migliavacca, M.,
768 Arain, M. A., and Janssens, I. A.: Evaluating the convergence between eddy-covariance and biometric methods
769 for assessing carbon budgets of forests, *Nature Communications*, 7, 13717, 2016.

770 Camps-Valls, G., Verrelst, J., Munoz-Mari, J., Laparra, V., Mateo-Jimenez, F., and Gomez-Dans, J.: A Survey
771 on Gaussian Processes for Earth-Observation Data Analysis: A Comprehensive Investigation, *IEEE Geoscience
772 and Remote Sensing Magazine*, 4, 58-78, 2016.

773 Chen, J. M., Mo, G., Pisek, J., Liu, J., Deng, F., Ishizawa, M., and Chan, D.: Effects of foliage clumping on the
774 estimation of global terrestrial gross primary productivity, *Global Biogeochemical Cycles*, 26, 2012.

775 Chevallier, F., Fisher, M., Peylin, P., Serrar, S., Bousquet, P., Bréon, F. M., Chédin, A., and Ciais, P.: Inferring
776 CO₂ sources and sinks from satellite observations: Method and application to TOVS data, *Journal of
777 Geophysical Research: Atmospheres*, 110, 2005.

778 Chevallier, F., Remaud, M., O'Dell, C. W., Baker, D., Peylin, P., and Cozic, A.: Objective evaluation of surface-
779 and satellite-driven CO₂ atmospheric inversions, *Atmos. Chem. Phys. Discuss.*, 2019, 1-28, 2019.

780 Damianou, A. and Lawrence, N.: Deep Gaussian Processes, Scottsdale, AZ, USA2013, 207-215.

781 Doelling, D. R., Loeb, N. G., Keyes, D. F., Nordeen, M. L., Morstad, D., Nguyen, C., Wielicki, B. A., Young, D.
782 F., and Sun, M.: Geostationary Enhanced Temporal Interpolation for CERES Flux Products, *J. Atmos. Ocean.
783 Technol.*, 30, 1072-1090, 2013.

784 Flach, M., Sippel, S., Gans, F., Bastos, A., Brenning, A., Reichstein, M., and Mahecha, M. D.: Contrasting
785 biosphere responses to hydrometeorological extremes: revisiting the 2010 western Russian heatwave,
786 *Biogeosciences*, 15, 6067-6085, 2018.

787 Frankenberg, C., Fisher, J. B., Worden, J., Badgley, G., Saatchi, S. S., Lee, J.-E., Toon, G. C., Butz, A., Jung,
788 M., Kuze, A., and Yokota, T.: New global observations of the terrestrial carbon cycle from GOSAT: Patterns of
789 plant fluorescence with gross primary productivity, *Geophysical Research Letters*, 38, L17706,
790 doi:17710.11029/12011GL048738, 2011.

791 Friedl, M. A., Sulla-Menashe, D., Tan, B., Schneider, A., Ramankutty, N., Sibley, A., and Huang, X.: MODIS
792 Collection 5 global land cover: Algorithm refinements and characterization of new datasets, *Remote Sensing of
793 Environment*, 114, 168-182, 2010.

794 Fu, Z., Gerken, T., Bromley, G., Araújo, A., Bonal, D., Burban, B., Ficklin, D., Fuentes, J. D., Goulden, M.,
795 Hirano, T., Kosugi, Y., Liddell, M., Nicolini, G., Niu, S., Rouspard, O., Stefani, P., Mi, C., Tofte, Z., Xiao, J.,
796 Valentini, R., Wolf, S., and Stoy, P. C.: The surface-atmosphere exchange of carbon dioxide in tropical
797 rainforests: Sensitivity to environmental drivers and flux measurement methodology, *Agricultural and Forest
798 Meteorology*, 263, 292-307, 2018.

799 Gaubert, B., Stephens, B. B., Basu, S., Chevallier, F., Deng, F., Kort, E. A., Patra, P. K., Peters, W., Rödenbeck,
800 C., Saeki, T., Schimel, D., Van der Laan-Luijkx, I., Wofsy, S., and Yin, Y.: Global atmospheric CO₂ inverse
801 models converging on neutral tropical land exchange, but disagreeing on fossil fuel and atmospheric growth rate,
802 *Biogeosciences*, 16, 117-134, 2019.

803 Guanter, L., Zhang, Y. G., Jung, M., Joiner, J., Voigt, M., Berry, J. A., Frankenberg, C., Huete, A. R., Zarco-
804 Tajada, P., Lee, J.-E., Moran, M. S., Ponce-Campos, G., Beer, C., Camps-Valls, G., Buchmann, N., Gianelle, D.,
805 Klumpp, K., Cescatti, A., Baker, J. M., and Griffis, T. J.: Global and time-resolved monitoring of crop
806 photosynthesis with chlorophyll fluorescence, *Proceedings of the National Academy of Sciences of the United
807 States of America*, 111, E1327-E1333, 2014.

808 Harris, I. C.: CRU JRA v1.1: A forcings dataset of gridded land surface blend of Climatic Research Unit (CRU)
809 and Japanese reanalysis (JRA) data. Unit, U. o. E. A. C. R. (Ed.), 2019.

810 Hayek, M. N., Wehr, R., Longo, M., Hutrya, L. R., Wiedemann, K., Munger, J. W., Bonal, D., Saleska, S. R.,
811 Fitzjarrald, D. R., and Wofsy, S. C.: A novel correction for biases in forest eddy covariance carbon balance,
812 *Agricultural and Forest Meteorology*, 250-251, 90-101, 2018.

813 Hellevang, H. and Aagaard, P.: Constraints on natural global atmospheric CO₂ fluxes from 1860 to 2010 using a
814 simplified explicit forward model, *Scientific Reports*, 5, 17352, 2015.

815 Huffman, G. J., Adler, R. F., Morrissey, M. M., Bolvin, D. T., Curtis, S., Joyce, R., McGavock, B., and
816 Susskind, J.: Global Precipitation at One-Degree Daily Resolution from Multisatellite Observations, *Journal of*
817 *Hydrometeorology*, 2, 36-50, 2001.

818 Humphrey, V., Zscheischler, J., Ciais, P., Gudmundsson, L., Sitch, S., and Seneviratne, S. I.: Sensitivity of
819 atmospheric CO₂ growth rate to observed changes in terrestrial water storage, *Nature*, 560, 628-631, 2018.

820 Hutyra, L. R., Munger, J. W., Hammond-Pyle, E., Saleska, S. R., Restrepo-Coupe, N., Daube, B. C., de
821 Camargo, P. B., and Wofsy, S. C.: Resolving systematic errors in estimates of net ecosystem exchange of CO₂
822 and ecosystem respiration in a tropical forest biome, *Agricultural and Forest Meteorology*, 148, 1266-1279,
823 2008.

824 Jiang, C. and Ryu, Y.: Multi-scale evaluation of global gross primary productivity and evapotranspiration
825 products derived from Breathing Earth System Simulator (BESS), *Remote Sensing of Environment*, 186, 528-
826 547, 2016.

827 Joiner, J., Yoshida, Y., Zhang, Y., Duveiller, G., Jung, M., Lyapustin, A., Wang, Y., and Tucker, J. C.:
828 Estimation of Terrestrial Global Gross Primary Production (GPP) with Satellite Data-Driven Models and Eddy
829 Covariance Flux Data, *Remote Sensing*, 10, 2018.

830 Jung, M., Koirala, S., Weber, U., Ichii, K., Gans, F., Camps-Valls, G., Papale, D., Schwalm, C., Tramontana, G.,
831 and Reichstein, M.: The FLUXCOM ensemble of global land-atmosphere energy fluxes, *Scientific Data*, in
832 press, 2019.

833 Jung, M., Reichstein, M., and Bondeau, A.: Towards global empirical upscaling of FLUXNET eddy covariance
834 observations: validation of a model tree ensemble approach using a biosphere model, *Biogeosciences*, 6, 2001-
835 2013, 2009.

836 Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneth, A., Bernhofer,
837 C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law, B. E., Lindroth, A.,
838 Merbold, L., Montagnani, L., Moors, E. J., Papale, D., Sottocornola, M., Vaccari, F., and Williams, C.: Global
839 patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible heat derived from eddy
840 covariance, satellite, and meteorological observations, *Journal of Geophysical Research - Biogeosciences*, 116,
841 G00J07, doi:10.1029/2010JG001566, 2011.

842 Jung, M., Reichstein, M., Schwalm, C. R., Huntingford, C., Sitch, S., Ahlström, A., Arneth, A., Camps-Valls, G.,
843 Ciais, P., Friedlingstein, P., Gans, F., Ichii, K., Jain, A. K., Kato, E., Papale, D., Poulter, B., Raduly, B.,
844 Rödenbeck, C., Tramontana, G., Viovy, N., Wang, Y.-P., Weber, U., Zaehle, S., and Zeng, N.: Compensatory
845 water effects link yearly global land CO₂ sink changes to temperature, *Nature*, 541, 516-520, 2017.

846 Keenan, T. F., Migliavacca, M., Papale, D., Baldocchi, D., Reichstein, M., Torn, M., and Wutzler, T.:
847 Widespread inhibition of daytime ecosystem respiration, *Nature Ecology & Evolution*, 3, 407-415, 2019.

848 Kim, H.: Global Soil Wetness Project Phase 3 Atmospheric Boundary Conditions (Experiment 1) [Data set].
849 (DIAS), D. I. a. A. S. (Ed.), 2017.

850 Kirschbaum, M. U. F., Zeng, G., Ximenes, F., Giltrap, D. L., and Zeldis, J. R.: Towards a more complete
851 quantification of the global carbon cycle, *Biogeosciences*, 16, 831-846, 2019.

852 Koffi, E. N., Rayner, P. J., Scholze, M., and Beer, C.: Atmospheric constraints on gross primary productivity and
853 net ecosystem productivity: Results from a carbon-cycle data assimilation system, *Global Biogeochemical*
854 *Cycles*, 26, 2012.

855 Köhler, P., Guanter, L., and Joiner, J.: A linear method for the retrieval of sun-induced chlorophyll fluorescence
856 from GOME-2 and SCIAMACHY data, *Atmos. Meas. Tech.*, 8, 2589-2608, 2015.

857 Lasslop, G., Reichstein, M., Papale, D., Richardson, A. D., Arneth, A., Barr, A., Stoy, P., and Wohlfahrt, G.:
858 Separation of net ecosystem exchange into assimilation and respiration using a light response curve approach:
859 critical issues and global evaluation, *Global Change Biology*, 16, 187-208, 2010.

860 Le Quéré, C., Andrew, R. M., Friedlingstein, P., Sitch, S., Hauck, J., Pongratz, J., Pickers, P. A., Korsbakken, J.
861 I., Peters, G. P., Canadell, J. G., Arneth, A., Arora, V. K., Barbero, L., Bastos, A., Bopp, L., Chevallier, F.,
862 Chini, L. P., Ciais, P., Doney, S. C., Gkritzalis, T., Goll, D. S., Harris, I., Haverd, V., Hoffman, F. M., Hoppema,
863 M., Houghton, R. A., Hurtt, G., Ilyina, T., Jain, A. K., Johannessen, T., Jones, C. D., Kato, E., Keeling, R. F.,
864 Goldewijk, K. K., Landschützer, P., Lefèvre, N., Lienert, S., Liu, Z., Lombardozzi, D., Metzl, N., Munro, D. R.,
865 Nabel, J. E. M. S., Nakaoka, S. I., Neill, C., Olsen, A., Ono, T., Patra, P., Peregon, A., Peters, W., Peylin, P.,
866 Pfeil, B., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson, E., Rocher, M., Rödenbeck, C., Schuster,
867 U., Schwinger, J., Séférian, R., Skjelvan, I., Steinhoff, T., Sutton, A., Tans, P. P., Tian, H., Tilbrook, B.,
868 Tubiello, F. N., van der Laan-Luijkx, I. T., van der Werf, G. R., Viovy, N., Walker, A. P., Wiltshire, A. J.,
869 Wright, R., Zaehle, S., and Zheng, B.: Global Carbon Budget 2018, *Earth Syst. Sci. Data*, 10, 2141-2194, 2018.

870 LeCun, Y., Bengio, Y., and Hinton, G.: Deep learning, *Nature*, 521, 436, 2015.

871 Leuning, R., van Gorsel, E., Massman, W. J., and Isaac, P. R.: Reflections on the surface energy imbalance
872 problem, *Agricultural and Forest Meteorology*, 156, 65-74, 2012.

873 Liang, M.-C., Mahata, S., Laskar, A. H., Thiemens, M. H., and Newman, S.: Oxygen isotope anomaly in
874 tropospheric CO₂ and implications for CO₂ residence time in the atmosphere and gross primary productivity,
875 *Scientific Reports*, 7, 13180, 2017.

876 MacBean, N., Maignan, F., Bacour, C., Lewis, P., Peylin, P., Guanter, L., Köhler, P., Gómez-Dans, J., and
877 Disney, M.: Strong constraint on modelled global carbon uptake using solar-induced chlorophyll fluorescence
878 data, *Scientific Reports*, 8, 1973, 2018.

879 Marcolla, B., Rödenbeck, C., and Cescatti, A.: Patterns and controls of inter-annual variability in
880 the terrestrial carbon budget, *Biogeosciences*, 14, 3815-3829, 2017.

881 McHugh, I. D., Beringer, J., Cunningham, S. C., Baker, P. J., Cavagnaro, T. R., Mac Nally, R., and Thompson,
882 R. M.: Interactions between nocturnal turbulent flux, storage and advection at an “ideal” eucalypt woodland site,
883 *Biogeosciences*, 14, 3027-3050, 2017.

884 Migliavacca, M., Perez-Priego, O., Rossini, M., El-Madany, T. S., Moreno, G., van der Tol, C., Rascher, U.,
885 Berninger, A., Bessenbacher, V., Burkart, A., Carrara, A., Fava, F., Guan, J.-H., Hammer, T. W., Henkel, K.,
886 Juarez-Alcalde, E., Julitta, T., Kolle, O., Martín, M. P., Musavi, T., Pacheco-Labrador, J., Pérez-Burgueño, A.,
887 Wutzler, T., Zaehle, S., and Reichstein, M.: Plant functional traits and canopy structure control the relationship
888 between photosynthetic CO₂ uptake and far-red sun-induced fluorescence in a Mediterranean grassland under
889 different nutrient availability, *New Phytologist*, 214, 1078-1091, 2017.

890 Morales, P., Sykes, M. T., Prentice, I. C., Smith, P., Smith, B., Bugmann, H., Zierl, B., Friedlingstein, P., Viovy,
891 N., Sabaté, S., Sánchez, A., Pla, E., Gracia, C. A., Sitch, S., Arneth, A., and Ogee, J.: Comparing and evaluating
892 process-based ecosystem model predictions of carbon and water fluxes in major European forest biomes, *Global
893 Change Biology*, 11, 2211-2233, 2005.

894 Musavi, T., Migliavacca, M., Reichstein, M., Kattge, J., Wirth, C., Black, T. A., Janssens, I., Knohl, A., Loustau,
895 D., Rouspard, O., Varlagin, A., Rambal, S., Cescatti, A., Gianelle, D., Kondo, H., Tamrakar, R., and Mahecha,
896 M. D.: Stand age and species richness dampen interannual variation of ecosystem-level photosynthetic capacity,
897 *Nature Ecology & Evolution*, 1, 0048, 2017.

898 Norton, A. J., Rayner, P. J., Koffi, E. N., Scholze, M., Silver, J. D., and Wang, Y. P.: Estimating global gross
899 primary productivity using chlorophyll fluorescence and a data assimilation system with the BETHY-SCOPE
900 model, *Biogeosciences Discuss.*, 2019, 1-45, 2019.

901 Orth, R., Destouni, G., Jung, M., and Reichstein, M.: Large-scale biospheric drought response intensifies linearly
902 with drought duration, *Biogeosciences Discuss.*, 2019, 1-25, 2019.

903 Pan, S. J. and Yang, Q.: A Survey on Transfer Learning, *IEEE Transactions on Knowledge and Data
904 Engineering*, 22, 1345-1359, 2010.

905 Pan, Y., Birdsey, R. A., Fang, J., Houghton, R., Kauppi, P. E., Kurz, W. A., Phillips, O. L., Shvidenko, A.,
906 Lewis, S. L., Canadell, J. G., Ciais, P., Jackson, R. B., Pacala, S. W., McGuire, A. D., Piao, S., Rautiainen, A.,
907 Sitch, S., and Hayes, D.: A Large and Persistent Carbon Sink in the World’s Forests, *Science*, 333, 988, 2011.

908 Papale, D., Black, T. A., Carvalhais, N., Cescatti, A., Chen, J., Jung, M., Kiely, G., Lasslop, G., Mahecha, M. D.,
909 Margolis, H., Merbold, L., Montagnani, L., Moors, E., Olesen, J. E., Reichstein, M., Tramontana, G., van Gorsel,
910 E., Wohlfahrt, G., and Ráduly, B.: Effect of spatial sampling from European flux towers for estimating carbon
911 and water fluxes with artificial neural networks, *Journal of Geophysical Research: Biogeosciences*, 120, 1941-
912 1957, 2015.

913 Papale, D., Reichstein, M., Aubinet, M., Canfora, E., Bernhofer, C., Kutsch, W., Longdoz, B., Rambal, S.,
914 Valentini, R., Vesala, T., and Yakir, D.: Towards a standardized processing of Net Ecosystem Exchange
915 measured with eddy covariance technique: algorithms and uncertainty estimation, *Biogeosciences*, 3, 571-583,
916 2006.

917 Peters, W., Krol, M. C., Van Der Werf, G. R., Houweling, S., Jones, C. D., Hughes, J., Schaefer, K., Masarie, K.
918 A., Jacobson, A. R., Miller, J. B., Cho, C. H., Ramonet, M., Schmidt, M., Ciattaglia, L., Apadula, F., Heltai, D.,
919 Meinhardt, F., Di Sarra, A. G., Piacentino, S., Sferlazzo, D., Aalto, T., Hatakka, J., Ström, J., Haszpra, L.,
920 Meijer, H. A. J., Van Der Laan, S., Neubert, R. E. M., Jordan, A., Rodó, X., Morguá, J. A., Vermeulen, A. T.,
921 Popa, E., Rozanski, K., Zimnoch, M., Manning, A. C., Leuenberger, M., Uglietti, C., Dolman, A. J., Ciais, P.,
922 Heimann, M., and Tans, P. P.: Seven years of recent European net terrestrial carbon dioxide exchange
923 constrained by atmospheric observations, *Global Change Biology*, 16, 1317-1337, 2010.

924 Peylin, P., Law, R. M., Gurney, K. R., Chevallier, F., Jacobson, A. R., Maki, T., Niwa, Y., Patra, P. K., Peters,
925 W., Rayner, P. J., Roedenbeck, C., van der Laan-Luijkx, I. T., and Zhang, X.: Global atmospheric carbon
926 budget: results from an ensemble of atmospheric CO₂ inversions, *Biogeosciences*, 10, 6699-6720, 2013.

927 Porcar-Castell, A., Tyystjärvi, E., Atherton, J., van der Tol, C., Flexas, J., Pfündel, E. E., Moreno, J.,
928 Frankenberg, C., and Berry, J. A.: Linking chlorophyll fluorescence to photosynthesis for remote sensing
929 applications: mechanisms and challenges, *Journal of Experimental Botany*, 65, 4065-4095, 2014.

930 Pugh, T. A. M., Arneth, A., Kautz, M., Poulter, B., and Smith, B.: Important role of forest disturbances in the
931 global biomass turnover and carbon sinks, *Nature Geoscience*, 12, 730-735, 2019.

932 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat: Deep learning
933 and process understanding for data-driven Earth system science, *Nature*, 566, 195-204, 2019.

934 Reichstein, M., Falge, E., Baldocchi, D., Papale, D., Valentini, R., Aubinet, M., Berbigier, P., Bernhofer, C.,
935 Buchmann, N., Gilmanov, T., Granier, A., Grünwald, T., Havránková, K., Janous, D., Knohl, A., Laurela, T.,
936 Lohila, A., Loustau, D., Matteucci, G., Meyers, T., Miglietta, F., Ourcival, J.-M., Rambal, S., Rotenberg, E.,

937 Sanz, M., Seufert, G., Vaccari, F., Vesala, T., and Yakir, D.: On the separation of net ecosystem exchange into
938 assimilation and ecosystem respiration: review and improved algorithm, *Global Change Biology*, 11, 1424-1439,
939 2005.

940 Rödenbeck, C., Zaehle, S., Keeling, R., and Heimann, M.: How does the terrestrial carbon exchange respond to
941 inter-annual climatic variations? A quantification based on atmospheric CO₂ data, *Biogeosciences*, 15, 2481-
942 2498, 2018.

943 Saleska, S. R., Miller, S. D., Matross, D. M., Goulden, M. L., Wofsy, S. C., da Rocha, H. R., de Camargo, P. B.,
944 Crill, P., Daube, B. C., de Freitas, H. C., Hutyra, L., Keller, M., Kirchhoff, V., Menton, M., Munger, J. W., Pyle,
945 E. H., Rice, A. H., and Silva, H.: Carbon in Amazon Forests: Unexpected Seasonal Fluxes and Disturbance-
946 Induced Losses, *Science*, 302, 1554, 2003.

947 Schimel, D., Pavlick, R., Fisher, J. B., Asner, G. P., Saatchi, S., Townsend, P., Miller, C., Frankenberg, C.,
948 Hibbard, K., and Cox, P.: Observing terrestrial ecosystems and the carbon cycle from space, *Global Change
949 Biology*, 21, 1762-1776, 2015.

950 Sitch, S., Friedlingstein, P., Gruber, N., Jones, S. D., Murray-Tortarolo, G., Ahlstrom, A., Doney, S. C., Graven,
951 H., Heinze, C., Huntingford, C., Levis, S., Levy, P. E., Lomas, M., Poulter, B., Viovy, N., Zaehle, S., Zeng, N.,
952 Arneeth, A., Bonan, G., Bopp, L., Canadell, J. G., Chevallier, F., Ciais, P., Ellis, R., Gloor, M., Peylin, P., Piao, S.
953 L., Le Quere, C., Smith, B., Zhu, Z., and Myneni, R.: Recent trends and drivers of regional sources and sinks of
954 carbon dioxide, *Biogeosciences*, 12, 653-679, 2015.

955 Spielmann, F. M., Wohlfahrt, G., Hammerle, A., Kitz, F., Migliavacca, M., Alberti, G., Ibrom, A., El-Madany,
956 T. S., Gerdel, K., Moreno, G., Kolle, O., Karl, T., Peressotti, A., and Delle Vedove, G.: Gross Primary
957 Productivity of Four European Ecosystems Constrained by Joint CO₂ and COS Flux Measurements,
958 *Geophysical Research Letters*, 46, 5284-5293, 2019.

959 Stoy, P. C., Mauder, M., Foken, T., Marcolla, B., Boegh, E., Ibrom, A., Arain, M. A., Arneeth, A., Aurela, M.,
960 Bernhofer, C., Cescatti, A., Dellwik, E., Duce, P., Gianelle, D., van Gorsel, E., Kiely, G., Knohl, A., Margolis,
961 H., McCaughey, H., Merbold, L., Montagnani, L., Papale, D., Reichstein, M., Saunders, M., Serrano-Ortiz, P.,
962 Sottocornola, M., Spano, D., Vaccari, F., and Varlagin, A.: A data-driven analysis of energy balance closure
963 across FLUXNET research sites: The role of landscape scale heterogeneity, *Agricultural and Forest
964 Meteorology*, 171-172, 137-152, 2013.

965 Sun, Y., Frankenberg, C., Wood, J. D., Schimel, D. S., Jung, M., Guanter, L., Drewry, D. T., Verma, M., Porcar-
966 Castell, A., Griffis, T. J., Gu, L., Magney, T. S., Köhler, P., Evans, B., and Yuen, K.: OCO-2 advances
967 photosynthesis observation from space via solar-induced chlorophyll fluorescence, *Science*, 358, eaam5747,
968 2017.

969 Tamrakar, R., Rayment, M. B., Moyano, F., Mund, M., and Knohl, A.: Implications of structural diversity for
970 seasonal and annual carbon dioxide fluxes in two temperate deciduous forests, *Agricultural and Forest
971 Meteorology*, 263, 465-476, 2018.

972 Tramontana, G., Jung, M., Schwalm, C. R., Ichii, K., Camps-Valls, G., Ráduly, B., Reichstein, M., Arain, M. A.,
973 Cescatti, A., Kiely, G., Merbold, L., Serrano-Ortiz, P., Sickert, S., Wolf, S., and Papale, D.: Predicting carbon
974 dioxide and energy fluxes across global FLUXNET sites with regression algorithms, *Biogeosciences*, 13, 4291-
975 4313, 2016.

976 van der Laan-Luijkx, I. T., van der Velde, I. R., van der Veen, E., Tsuruta, A., Stanislawski, K.,
977 Babenhauserheide, A., Zhang, H. F., Liu, Y., He, W., Chen, H., Masarie, K. A., Krol, M. C., and Peters, W.: The
978 CarbonTracker Data Assimilation Shell (CTDAS) v1.0: implementation and global carbon balance 2001–2015,
979 *Geosci. Model Dev.*, 10, 2785-2800, 2017.

980 van Gorsel, E., Delpierre, N., Leuning, R., Black, A., Munger, J. W., Wofsy, S., Aubinet, M., Feigenwinter, C.,
981 Beringer, J., Bonal, D., Chen, B., Chen, J., Clement, R., Davis, K. J., Desai, A. R., Dragoni, D., Etzold, S.,
982 Grünwald, T., Gu, L., Heinesch, B., Hutyra, L. R., Jans, W. W. P., Kutsch, W., Law, B. E., Leclerc, M. Y.,
983 Mammarella, I., Montagnani, L., Noormets, A., Rebmann, C., and Wharton, S.: Estimating nocturnal ecosystem
984 respiration from the vertical turbulent flux and change in storage of CO₂, *Agricultural and Forest Meteorology*,
985 149, 1919-1930, 2009.

986 van Gorsel, E., Leuning, R., Cleugh, H. A., Keith, H., Kirschbaum, M. U. F., and Suni, T.: Application of an
987 alternative method to derive reliable estimates of nighttime respiration from eddy covariance measurements in
988 moderately complex topography, *Agricultural and Forest Meteorology*, 148, 1174-1180, 2008.

989 Walther, S., Duveiller, G., Jung, M., Guanter, L., Cescatti, A., and Camps-Valls, G.: Satellite Observations of the
990 Contrasting Response of Trees and Grasses to Variations in Water Availability, *Geophysical Research Letters*,
991 46, 1429-1440, 2019.

992 Walther, S., Voigt, M., Thum, T., Gonsamo, A., Zhang, Y. G., Kohler, P., Jung, M., Varlagin, A., and Guanter,
993 L.: Satellite chlorophyll fluorescence measurements reveal large-scale decoupling of photosynthesis and
994 greenness dynamics in boreal evergreen forests, *Global Change Biology*, 22, 2979-2996, 2016.

995 Wang, K. A., Pleiss, G., Gardner, J. R., Tyree, S., Weinberger, K. Q., and Wilson, A. G.: Exact Gaussian
996 Processes on a Million Data Points, arXiv, 2019. 2019.

997 Weedon, G. P., Balsamo, G., Bellouin, N., Gomes, S., Best, M. J., and Viterbo, P.: The WFDEI meteorological
998 forcing data set: WATCH Forcing Data methodology applied to ERA-Interim reanalysis data, *Water Resources*
999 *Research*, 50, 7505-7514, 2014.

1000 Wehr, R., Munger, J. W., McManus, J. B., Nelson, D. D., Zahniser, M. S., Davidson, E. A., Wofsy, S. C., and
1001 Saleska, S. R.: Seasonality of temperate forest photosynthesis and daytime respiration, *Nature*, 534, 680, 2016.

1002 Welp, L. R., Keeling, R. F., Meijer, H. A. J., Bollenbacher, A. F., Piper, S. C., Yoshimura, K., Francey, R. J.,
1003 Allison, C. E., and Wahlen, M.: Interannual variability in the oxygen isotopes of atmospheric CO₂ driven by El
1004 Niño, *Nature*, 477, 579, 2011.

1005 You, J., Li, X., Low, M., Lobell, D., and Ermon, S.: Deep Gaussian Process for Crop Yield Prediction Based on
1006 Remote Sensing Data, 2017.

1007 Yu, T., Sun, R., Xiao, Z., Zhang, Q., Liu, G., Cui, T., and Wang, J.: Estimation of Global Vegetation
1008 Productivity from Global LAnd Surface Satellite Data, *Remote Sensing*, 10, 2018.

1009 Yuan, W., Liu, S., Yu, G., Bonnefond, J.-M., Chen, J., Davis, K., Desai, A. R., Goldstein, A. H., Gianelle, D.,
1010 Rossi, F., Suyker, A. E., and Verma, S. B.: Global estimates of evapotranspiration and gross primary production
1011 based on MODIS and global meteorology data, *Remote Sensing of Environment*, 114, 1416-1431, 2010.

1012 Zhang, Y., Guanter, L., Berry, J. A., van der Tol, C., Yang, X., Tang, J., and Zhang, F.: Model-based analysis of
1013 the relationship between sun-induced chlorophyll fluorescence and gross primary production for remote sensing
1014 applications, *Remote Sensing of Environment*, 187, 145-155, 2016.

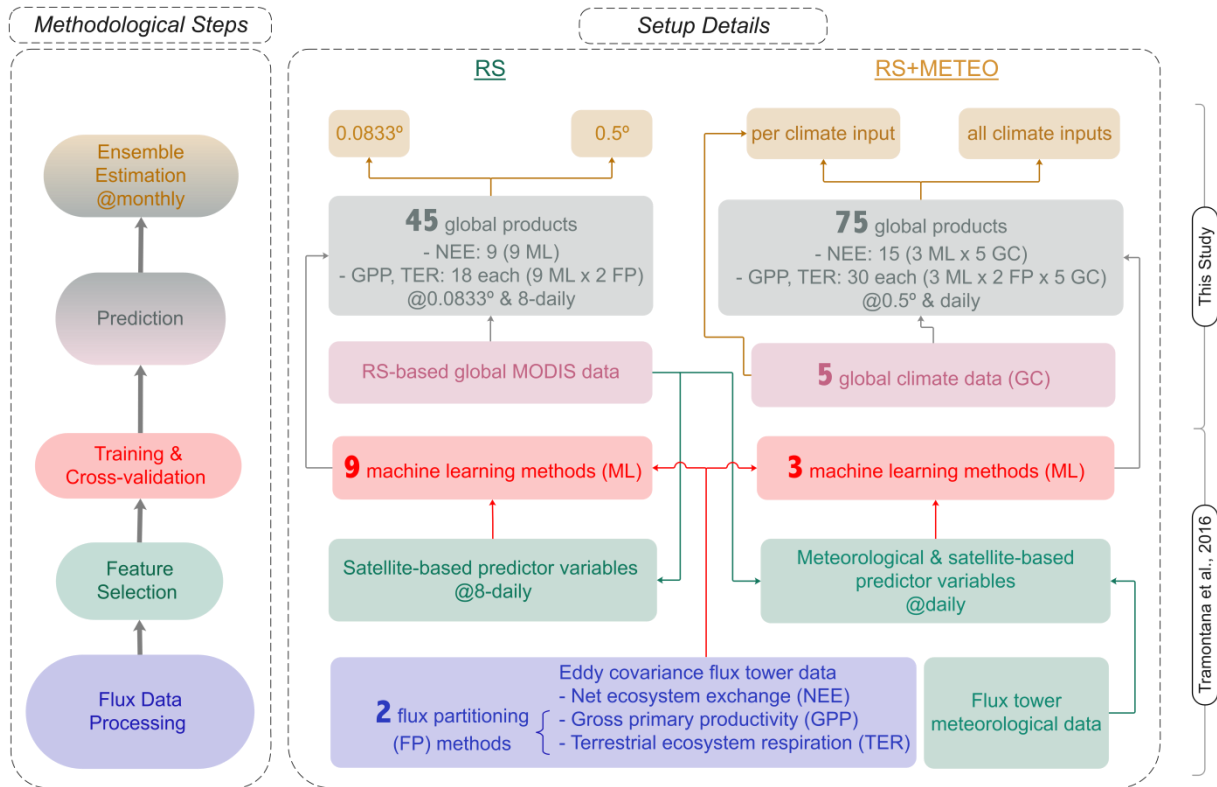
1015 Zhang, Y., Joiner, J., Gentile, P., and Zhou, S.: Reduced solar-induced chlorophyll fluorescence from GOME-2
1016 during Amazon drought caused by dataset artifacts, *Global Change Biology*, 24, 2229-2230, 2018.

1017 Zhao, M., Heinsch, F. A., Nemani, R. R., and Running, S. W.: Improvements of the MODIS terrestrial gross and
1018 net primary production global data set, *Remote Sensing of Environment*, 95, 164-176, 2005.

1019 Zscheischler, J., Mahecha, M. D., Avitabile, V., Calle, L., Carvalhais, N., Ciais, P., Gans, F., Gruber, N.,
1020 Hartmann, J., Herold, M., Ichii, K., Jung, M., Landschützer, P., Laruelle, G. G., Lauerwald, R., Papale, D.,
1021 Peylin, P., Poulter, B., Ray, D., Regnier, P., Rödenbeck, C., Roman-Cuesta, R. M., Schwalm, C., Tramontana,
1022 G., Tyukavina, A., Valentini, R., van der Werf, G., West, T. O., Wolf, J. E., and Reichstein, M.: Reviews and
1023 syntheses: An empirical spatiotemporal description of the global surface-atmosphere carbon fluxes:
1024 opportunities and data limitations, *Biogeosciences*, 14, 3685-3703, 2017.

1025

1027

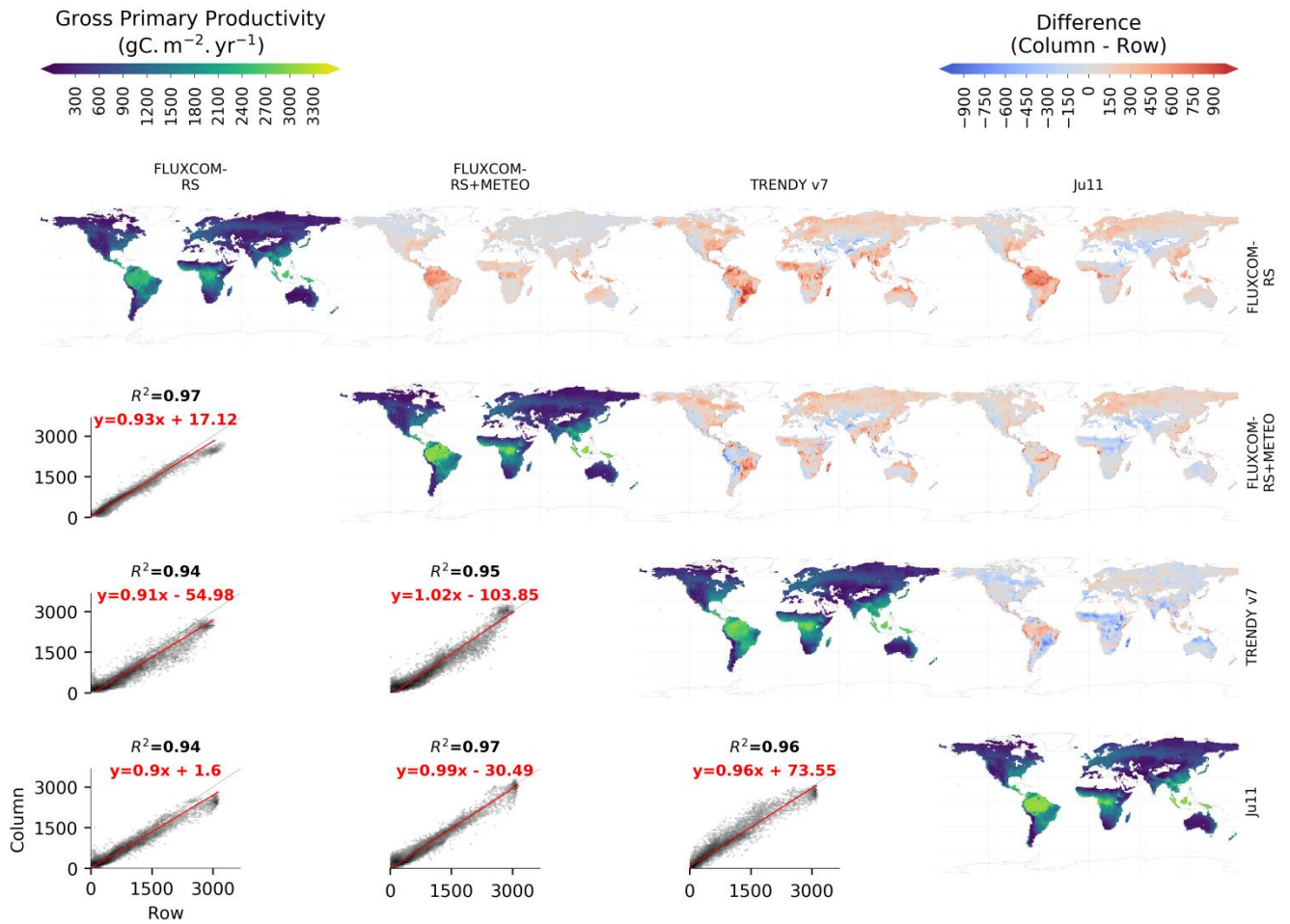


1028

1029

1030 **Figure 1: Schematic overview of the methodology and data products from the FLUXCOM initiative. The flow**
 1031 **diagram shows the methodological steps for the remote sensing -based (RS, left) and the remote sensing and**
 1032 **meteorological data -based (RS+METEO, right) FLUXCOM products. Final monthly ensemble products for NEE,**
 1033 **GPP, and TER from RS are available at 0.0833° and at 0.5° spatial resolution. Ensemble products from RS+METEO**
 1034 **are available per climate forcing (GC) data set as well as a pooled ensemble at 0.5° spatial resolution. All ensemble**
 1035 **products encompass ensemble members of different machine learning methods (ML, 9 for RS, 3 for RS+METEO) and**
 1036 **flux partitioning methods (FP, 2 for GPP and TER).**

1037

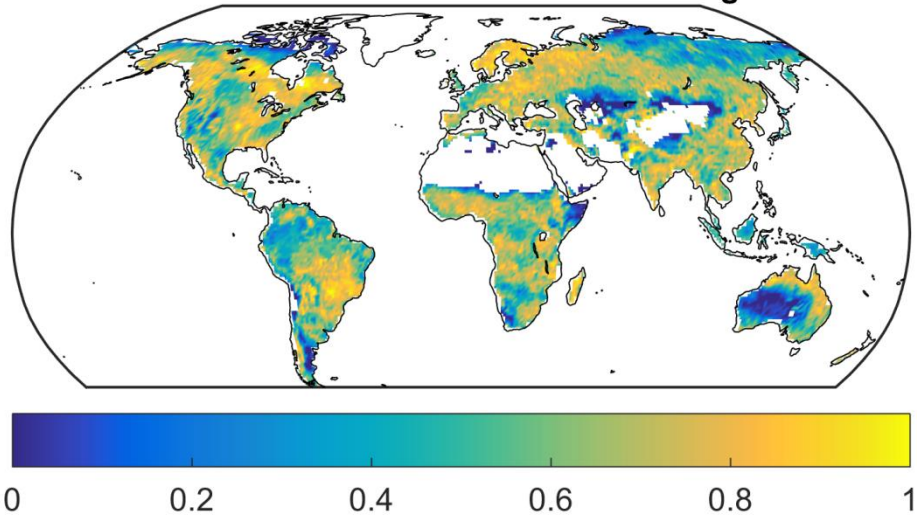


1038

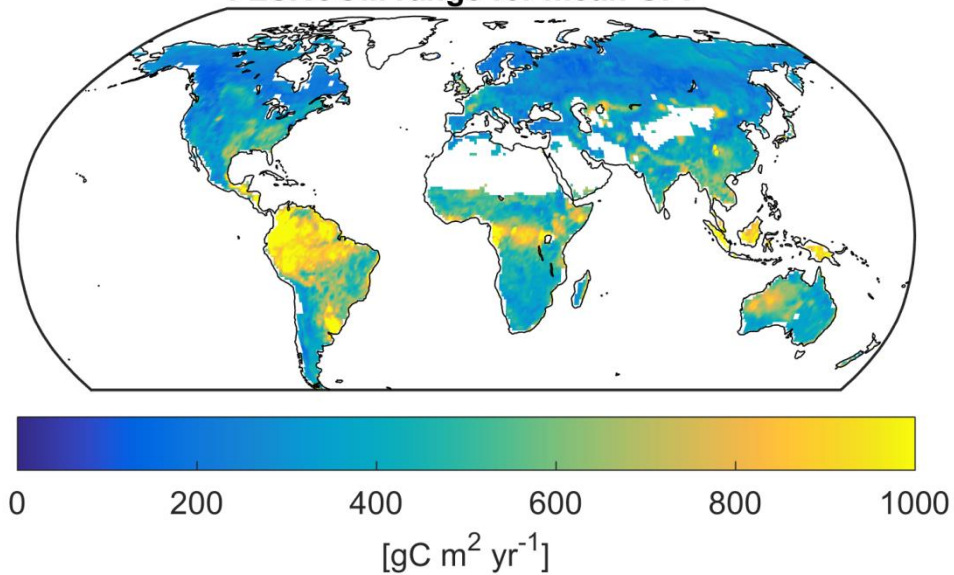
1039 **Figure 2: Comparisons of mean annual GPP at 1° spatial resolution for the period 2008-2010 of FLUXCOM ensemble**
 1040 **products with Ju11 and the mean of 16 TRENDY models. Diagonal: Maps of mean annual GPP. Above diagonal:**
 1041 **Maps of GPP differences (product along column – product along row). Below diagonal: 1:1 regression where the**
 1042 **shading shows point density. The red line and equations show the best fit line from total least square regression.**

1043

Fraction of TRENDY models outside FLUXCOM range for mean GPP

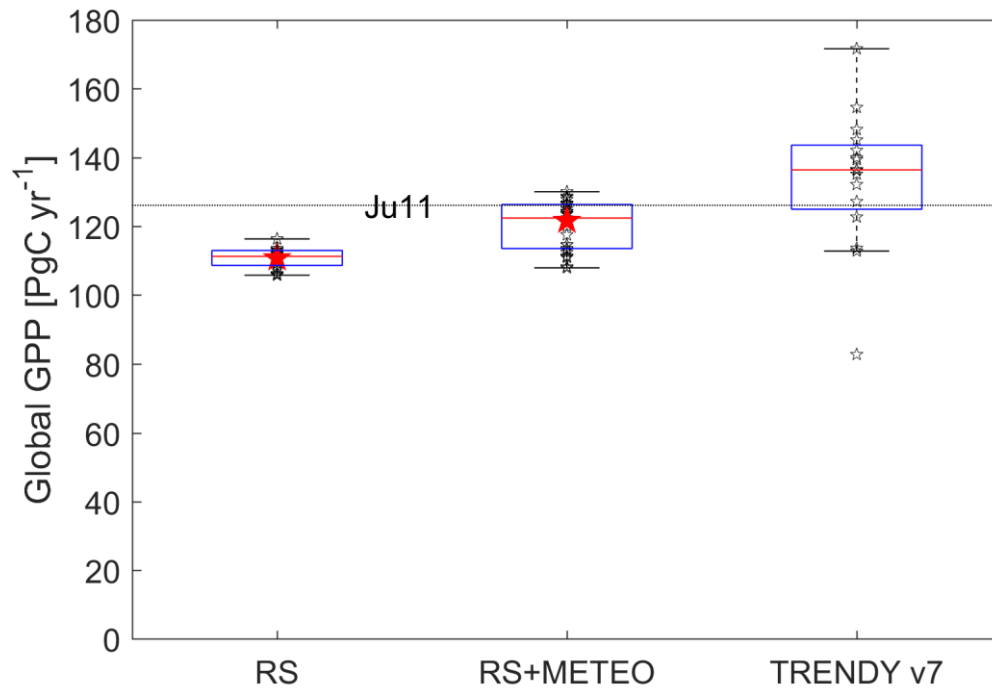


FLUXCOM range for mean GPP



1044

1045 **Figure 3: Map of the fraction of TRENDY models (n=16) with mean GPP outside the range of FLUXCOM estimates.**
1046 **The FLUXCOM range is calculated as the maximum minus minimum of all 48 FLUXCOM members from the union**
1047 **of the RS and RS+METEO members. Mean GPP was calculated for the period 2008-2010.**

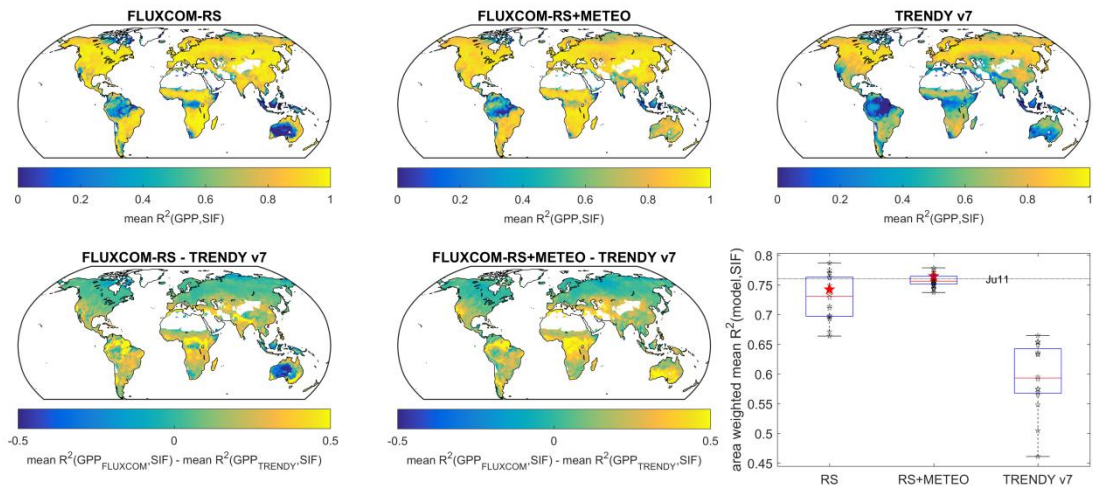


1048

1049 **Figure 4: Global GPP for FLUXCOM and TRENDY ensembles for the period 2008-2010. The box plots show the**
 1050 **median (red line), interquartile range (box) and total range (whiskers) of non-outliers (within median \pm 1.5**
 1051 **interquartile range) of individual ensemble members (open black stars). The filled red star presents the value of the**
 1052 **ensemble product (not available for TRENDY). The estimate of Ju11 is plotted as horizontal broken line.**

1053

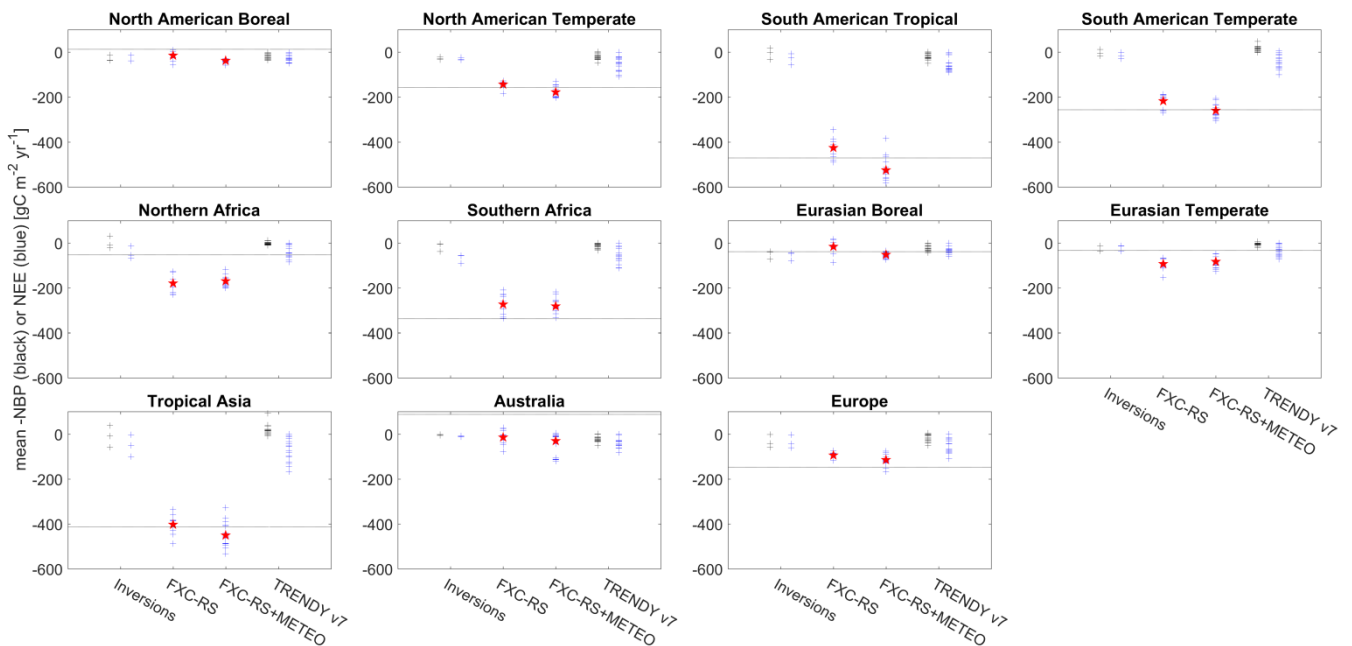
1054



1055

1056 **Figure 5: Consistency of seasonal GPP variations from FLUXCOM and TRENDY with SIF from GOME-2. Maps in the**
 1057 **top row show the mean R^2 between mean seasonal cycles for the period 2008-2010, averaged across all respective**
 1058 **ensemble members. Difference maps in the bottom row emphasize where FLUXCOM shows better (positive value)**
 1059 **and worse (negative value) consistency with SIF than TRENDY and are based on the maps in the top row. The**
 1060 **spatially averaged R^2 values for the different ensembles are summarized in the bottom right panel. The box plots**
 1061 **show the distribution of individual ensemble members (open black stars). The filled red star presents the value of the**
 1062 **ensemble product (not available for TRENDY). The estimate of Ju11 is plotted as horizontal broken line.**

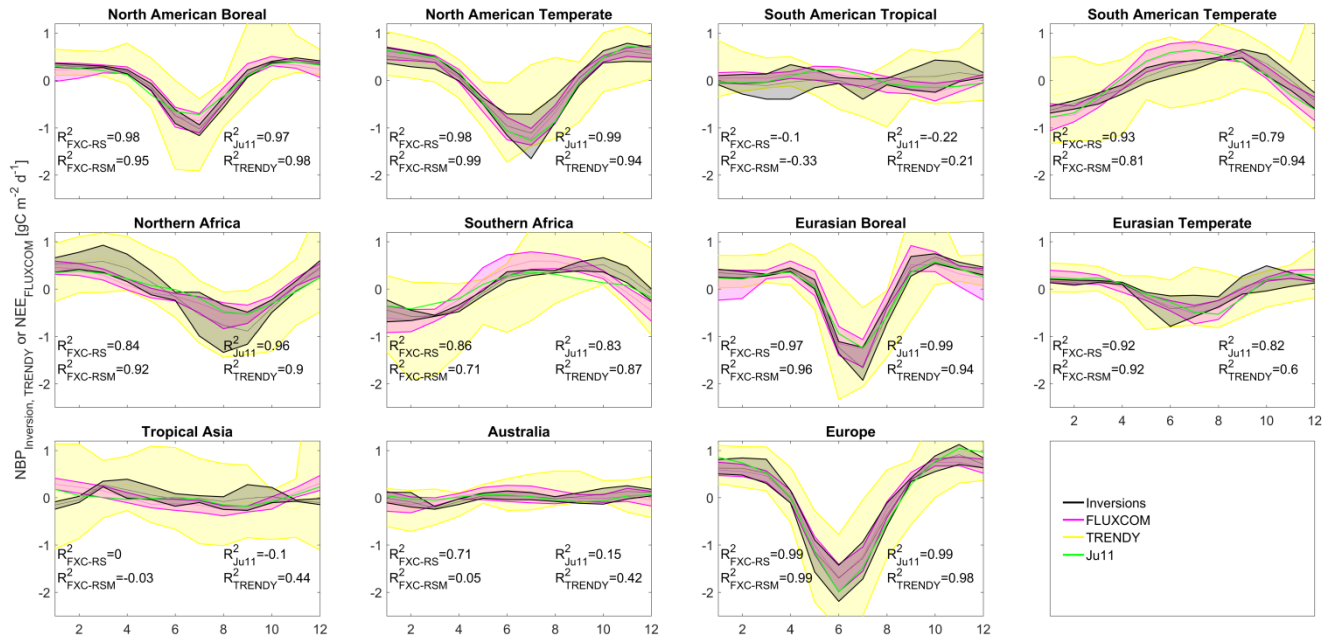
1063



1064

1065 **Figure 6: Mean annual net carbon release for the years 2008-2010 over TRANSKOM regions. Crosses refer to**
 1066 **individual ensemble members where a black colour refers to negative net biome productivity (NBP, not available for**
 1067 **FLUXCOM), and blue color refers to net ecosystem exchange (NEE). For inversions, NEE was approximated by**
 1068 **correcting NBP with fire emissions (see section 2.4.3). The filled red stars refer to estimates by the ensemble product**
 1069 **of FLUXCOM setups. The horizontal broken line indicates the estimate of Ju11.**

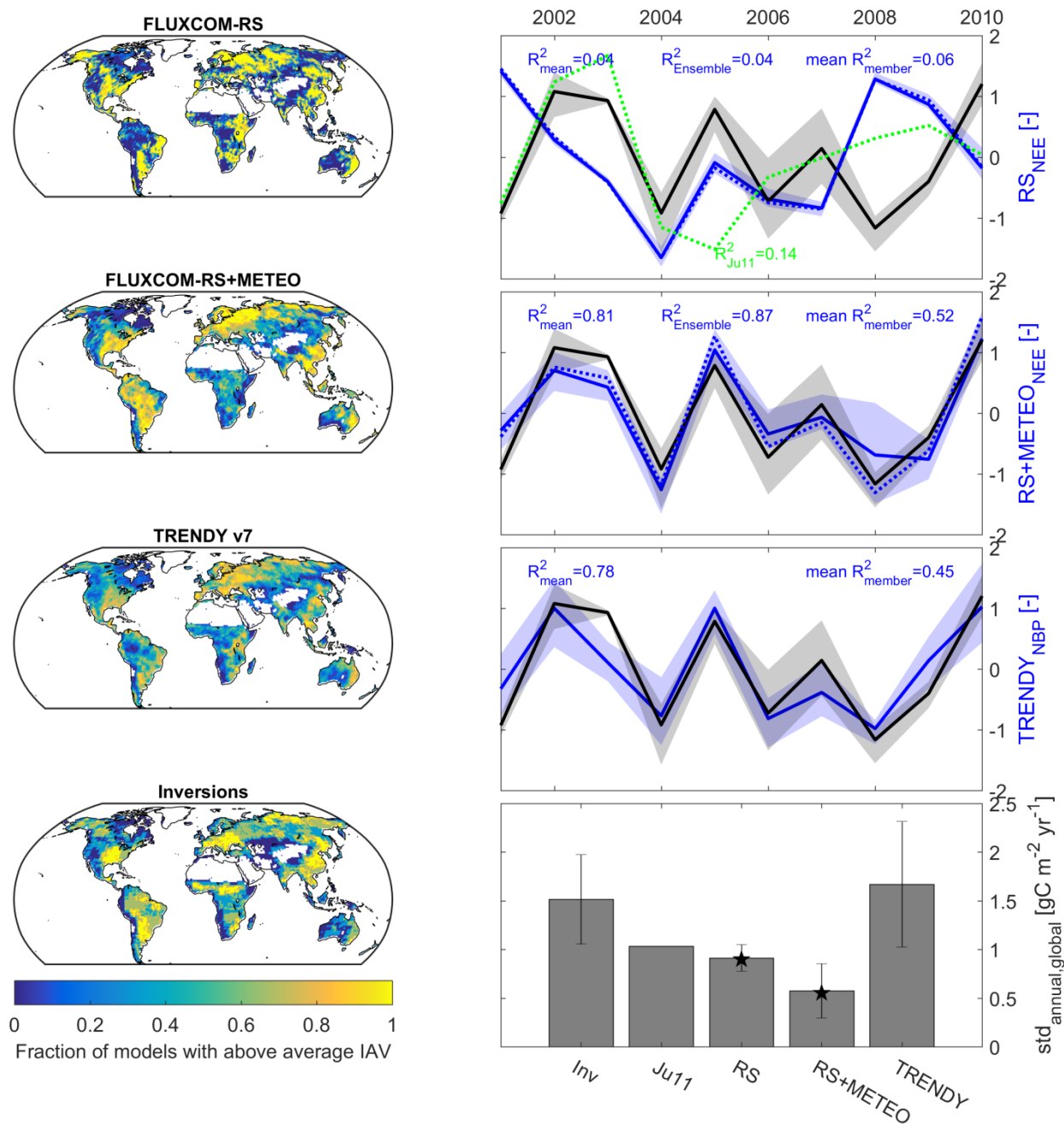
1070



1071

1072 **Figure 7: Mean seasonal variations of net land carbon release for the period 2008-2010 over TRANSCOM regions.**
 1073 **For inversions and TRENDY, -NBP was plotted, and for FLUXCOM, NEE was plotted. Please note that the region**
 1074 **specific mean was removed for each data set. Shading indicates the range of estimates (maximum – minimum). The**
 1075 **FLUXCOM range is based on the union of RS and RS+METEO ensemble members. R² values were calculated with**
 1076 **the mean of the inversions. The FLUXCOM RS and RS+METEO refer to the ensemble products (median), while that**
 1077 **for TRENDY refer to the model mean.**

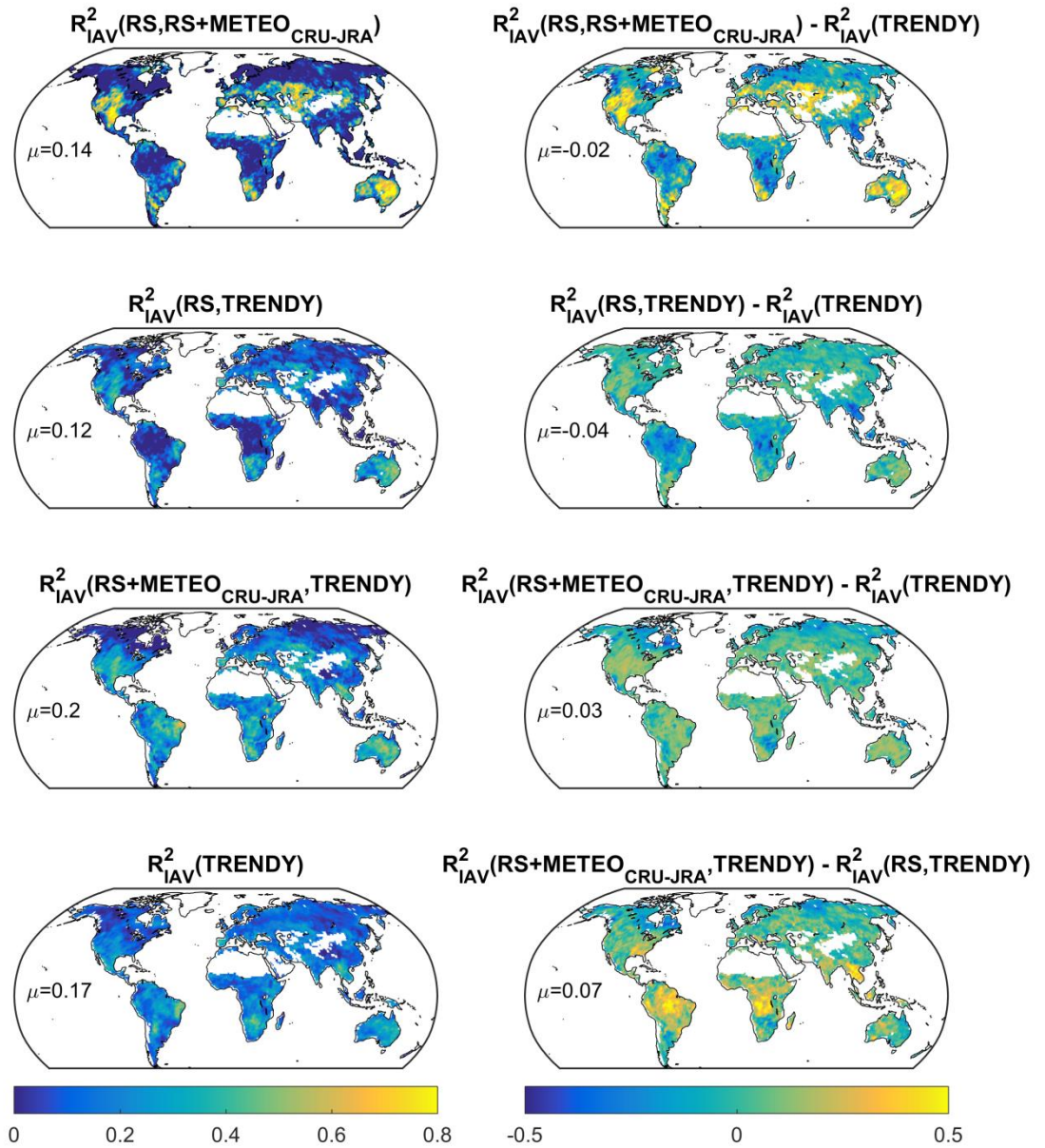
1078



1079

1080 **Figure 8: Interannual variability patterns of FLUXCOM NEE, TRENDY NBP, and NBP from three atmospheric**
 1081 **inversions for the period 2001-2010. Maps show the fraction of respective ensemble members with above average**
 1082 **interannual variability (standard deviation of annual values multiplied with land area). Time series plots show**
 1083 **detrended globally integrated annual NEE or NBP anomalies normalized by their standard deviation. The black line**
 1084 **is the mean of three inversions and the gray shading indicates their range. The blue solid lines are the means of the**
 1085 **considered ensembles; the blue dashed lines are the FLUXCOM ensemble products. R^2 values refer to the comparison**
 1086 **with the mean of inversions (black solid line). The bar chart in the bottom right panel shows the standard deviation of**
 1087 **detrended annual NEE or NBP for different data sets, averaged over the ensemble members and the error bar**
 1088 **indicates the standard deviation of the ensemble members. Black stars for FLUXCOM refer to the value for the**
 1089 **ensemble products.**

1090

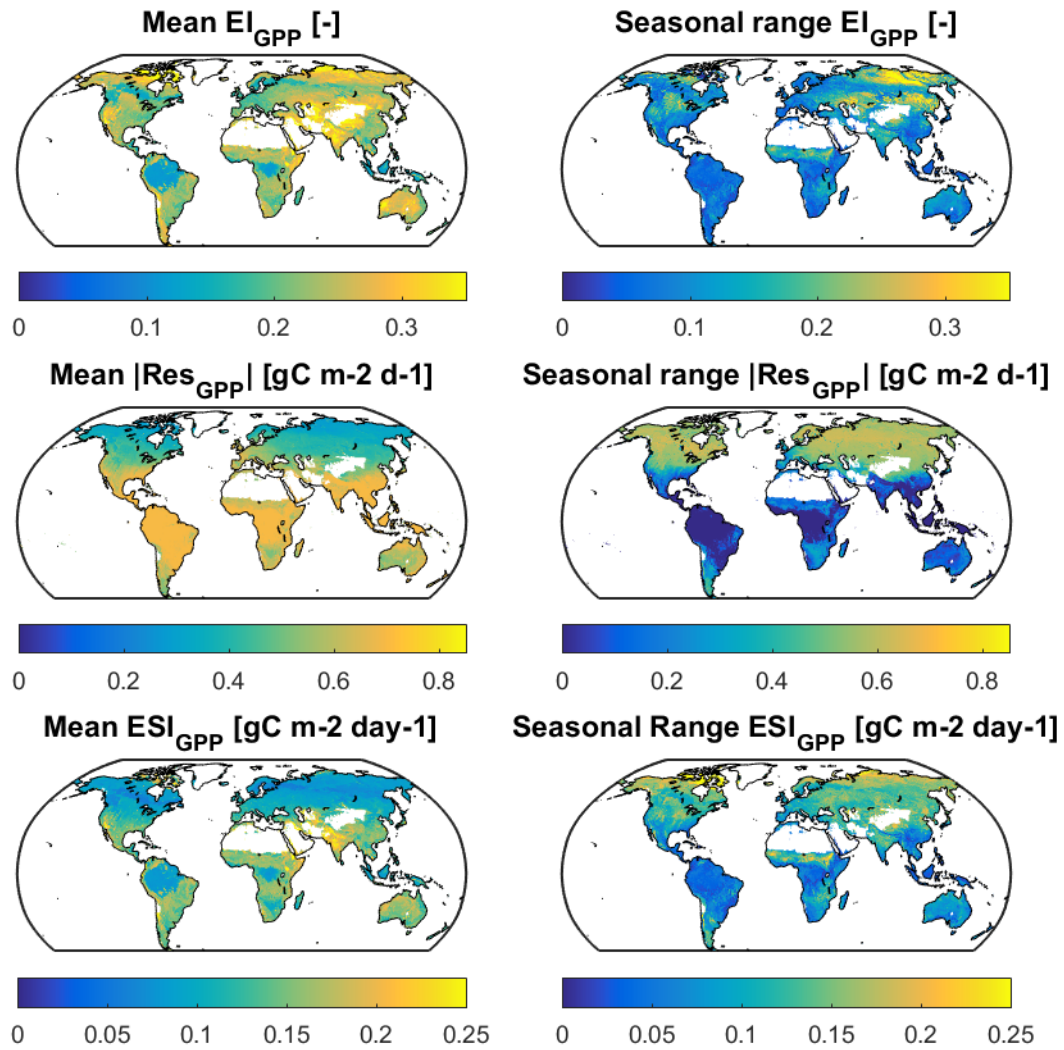


1091

1092 **Figure 9: Consistency between interannual variabilities (IAV) of local NEE from FLUXCOM setups and TRENDY**
 1093 **for the period 2001-2015.**

1094

1095



1096

1097 **Figure 10: Mean annual (2001-2015) and seasonal range (8-daily time step) of the Extrapolation Index (EI), the**
 1098 **expected mean absolute error of machine learning predictions, and the Extrapolation Severity Index (ESI, product of**
 1099 **the previous two) (see S2 for details) for GPP from FLUXCOM-RS.**

1100

1101

Meteorological forcing data set	Spatial Resolution	Temporal Coverage
CRU-JRA	0.5° x 0.5°	1950-2017
GSWP3	0.5° x 0.5°	1950-2010
WFDEI	0.5° x 0.5°	1979-2013
ERA-5	0.5° x 0.5°	1979-2018
CERES-GPCP	1.0° x 1.0° resampled to 0.5° x 0.5°	2001-2013

1102 **Table 1: Global meteorological forcing data sets used in FLUXCOM-RS+METEO.**

1103

1104

1105

1106

1107