# Ensemble modelling, uncertainty and robust predictions of organic carbon in long-term bare-fallow soils

*Model inter-comparison of soil organic carbon*

Farina, Roberta[1,*], Sándor, Renata[2,3], Abdalla, Mohamed[4], Álvaro-Fuentes, Jorge[5], Bechini, Luca[6], Bolinder, Martin A.[7], Brilli, Lorenzo[8], Chenu, Claire[9], Clivot, Hugues[10,11], De Antoni Migliorati, Massimiliano[12], Di Bene, Claudia[1], Dorich, Christopher D.[13], Ehrhardt, Fiona[14], Ferchaud, Fabien[10], Fitton, Nuala[4], Francaviglia, Rosa[1], Franko, Uwe[15], Giltrap, Donna.L.[16], Grant, Brian, B.[17], Guenet, Bertrand[18,19], Harrison, Matthew T.[20], Kirschbaum, Miko U.F.[16], Kuka, Katrin[21], Kulmala, Liisa[22], Liski, Jari[22], McGrath, Matthew J.[18], Meier, Elizabeth[23], Menichetti, Lorenzo[7], Moyano, Fernando[24], Nendel, Claas[25,29], Recous, Sylvie[26], Reibold, Nils[24], Shepherd, Anita[4,27] Smith, Ward N,[17], Smith, Pete[4], Soussana, Jean-François[14], Stella, Tommaso[25], Taghizadeh-Toosi, Arezoo.[28], Tsutskikh, Elena[25], Bellocchi, Gianni[3]

[1] CREA - Council for Agricultural Research and Economics, Research Centre for Agriculture and Environment, Rome, Italy

[2] Agricultural Institute, Centre for Agricultural Research, Martonvásár, Hungary

[3] Université Clermont Auvergne, INRAE, VetAgro Sup, UREP, Clermont-Ferrand, France

[4] University of Aberdeen, UK

[5] Spanish National Research Council (CSIC), Zaragoza, Spain

[6] Università degli Studi di Milano, Italy

[7] Swedish University of Agricultural Sciences, Uppsala, Sweden

[8] CNR-IBE, Institute of Bioeconomy, Florence, Italy

[9] Université Paris Saclay, INRAE, AgroParisTech, Paris, France

[10] INRAE, BioEcoAgro, F-02000, Barenton-Bugny, France

26    [11] Université de Lorraine, INRAE, LAE, F-68000, Colmar, France

27    [12] Queensland University of Technology, Brisbane, Australia

28    [13] Colorado State University, Fort Collins CO, USA

29    [14] INRAE, CODIR, 75007 Paris, France

30    [15] Helmholtz Centre for Environmental Research, Halle, Germany

31    [16] Manaaki Whenua - Landcare Research, Palmerston North, New Zealand

32    [17] Ottawa Research and Development Centre, Agriculture and Agri-Food, Ottawa, Canada

33    [18] Laboratoire des Sciences du Climat et de l'Environnement, LSCE/IPSL, CEA-CNRS-UVSQ,

34    Université Paris-Saclay, 91191 Gif-sur-Yvette, France

35    [19] Laboratoire de Géologie de l'ENS, PSL Research University, Paris, France

36    [20] Tasmanian Institute of Agriculture, Australia

37    [21] JKI - Federal Research Centre for Cultivated Plants, Braunschweig, Germany

38    [22] Finnish Meteorological Institute, Helsinki, Finland

39    [23] CSIRO, Brisbane, Australia

40    [24] University of Gottingen, Germany

41    [25] Leibniz Centre for Agricultural Landscape Research, Müncheberg, Germany,

42    [26] Université de Reims Champagne Ardenne, INRAE, FARE, Reims, France

43    [27] formerly Rothamsted Research, North Wyke, Devon, UK

44    [28] Department of Agroecology, Aarhus University, Tjele, Denmark

45    [29] University of Potsdam, Germany

46

47

48    *Corresponding author. Tel.: +39067005413; fax +39067005711

49    E-mail address: roberta.farina@crea.gov.it

50

51

**Abstract**

Simulation models represent soil organic carbon (SOC) dynamics in global carbon (C) cycle scenarios to support climate-change studies. It is imperative to increase confidence in long-term predictions of SOC dynamics by reducing the uncertainty in model estimates. To do this, we evaluated SOC simulated from an ensemble of 26 process-based C models by comparing simulations to experimental data from seven long-term bare-fallow (vegetation-free) plots at six sites in Denmark (two sites), France, Russia, Sweden and the United Kingdom. The decay of SOC in these plots has been monitored for decades since the last inputs of plant material, providing the opportunity to test decomposition without the continuous input of new organic material. The models were run independently over multi-year simulation periods (from 28 to 80 years) in a blind test with no calibration (Bln) and with three calibration scenarios, each providing different levels of information and/or allowing different levels of model fitting: a) calibrating decomposition parameters separately at each experimental site (Spe); b) using a generic, knowledge-based, parameterisation applicable in the Central European region (Gen); and c) using a combination of both a) and b) strategies (Mix). With this methodology, we addressed uncertainties from different modelling approaches with or without spin-up initialisation of SOC. Changes in the multi-model median (MMM) of SOC were used as descriptors of the ensemble performance. On average across sites, Gen proved adequate in describing changes in SOC, with MMM equal to average SOC (and standard deviation) of 39.2 ($\pm$15.5) Mg C ha$^{-1}$ compared to the observed mean of 36.0 ($\pm$19.7) Mg C ha$^{-1}$ (last observed year), indicating sufficiently reliable SOC estimates. This is important because moving to Mix (37.5$\pm$16.7 Mg C ha$^{-1}$) and Spe (36.8$\pm$19.8 Mg C ha$^{-1}$) provided only marginal gains in accuracy, but with these scenarios modellers would need to apply increasingly more knowledge and a greater calibration effort than in Gen, thereby limiting the wider applicability of models.

| LIST OF SYMBOLS AND ABBREVIATIONSSymbol/abbreviation | Long version | Explanation |
|---|---|---|
| *System variables* | | |
| C | Carbon | Chemical element with atomic number 6 |
| SOC | Soil organic carbon | Carbon stored in soil organic matter |
| SOM | Soil organic matter | The fraction of the soil that consists of plant, animal or microbial tissue in various stages of decomposition |
| N | Nitrogen | Chemical element with atomic number 7 |
| *Experimentation* | | |
| LTE | Long-term field experiment | Research facility providing data for monitoring trends and evaluating different agricultural management strategies over time |
| LTBF | Long-term bare-fallow experimental site | Research facility providing data for monitoring trends on bare-fallow soils |
| S1 | Site 1 | Askov (Denmark) – location 1 |
| S2 | Site 2 | Askov (Denmark) – location 2 |
| S3 | Site 3 | Grignon (France) |
| S4 | Site 4 | Kursk (Russia) |
| S5 | Site 5 | Rothamsted (United Kingdom) |
| S6 | Site 6 | Ultuna (Sweden) |
| S7 | Site 7 | Versailles (France) |
| *Modelling* | | |
| M01, …, M34 | Model 01, …, model 34 | Simulation models (M) anonymously coded from 1 to 34 |
| Bln | Blind | Uncalibrated simulations (blind test) |
| Gen | Generic | Generic simulation scenario |
| Mix | Mixed | Mixed simulation scenario |
| Spe | Specific | Specific simulation scenario |
| SP | Spin-up | Process of running the model from a set of conditions to initialise the state of C pools |
| NS | No spin-up | Any function (or analytical procedures) to make an initial partition of C pools (alternative to spin-up runs) |
| *Statistics* | | |

| | | |
|---|---|---|
| SD | Standard deviation | Variation amount of a set of data |
| MMM | Multi-model median | Median value of simulated data from different models |
| Obs | Observations | Observed data |
| RRMSE | Relative root mean square error | Aggregate magnitude of the errors in predictions relative to the mean of observations |
| EF | Modelling efficiency | Predictive power of a model with respect to the mean of observations |
| $R^2$ | Coefficient of determination | Proportion of the variance in the modelled data that is predictable from the observations |
| r | Pearson's correlation coefficient | Degree to which predictions and observations are linearly related |
| P(t) | Paired Student t-test probability of I-type error | Probability to reject the true null hypothesis of equal means of two samples of paired data (i.e. predictions and observations) |
| d | Index of agreement | Ratio of the mean square error and the potential error represented by the largest value that the squared difference of each prediction/observation pair can attain |
| $z$ | $z$-score transformation | Number of standard deviations by which the value of a raw score is above or below the mean value of the variable of interest |
| $sd$ | Standard deviation | Standard deviation units expressing $z$-scores |
| $sd_{obs}$ | Standard deviation of observations | Variation amount of a set of observed values |
| P | Predicted value | Value of a variable that is generated using a model |
| O | Observed value | Value of a variable that is actually observed |
| n | Number of predicted or observed values | Number of predicted/observed pairs |
| i | $i^{th}$ predicted or observed value | Subscript index of each predicted/observed pair |
| $\bar{O}$ | Mean of observed values | Arithmetic mean of actually observed data |

| | | |
|---|---|---|
| $\overline{P}$ | Mean of predicted values | Arithmetic mean of actually observed data |
| $\overline{D}$ | Mean difference | Arithmetic mean of the differences between predicted and observed values |
| $S_D$ | Standard deviation of the differences | Variation amount of a set of differences between predictions and observations |
| p | Probability of I-type error | Probability to reject the true null hypothesis of null correlation between two variable |

*Agro-climatic metrics*

| | | |
|---|---|---|
| Tamp | Temperature amplitude | Difference between the highest and the lowest temperature in a year |
| Tmax | Maximum air temperature | Average of the highest daily temperatures in a year |
| Prec | Precipitation | Annual precipitation total |
| $b^a$ | De Martonne-Gottman aridity index | Indicator of aridity including both annual and monthly temperature and precipitation |
| $hw^a$ | Heatwave frequency | Number of at least seven consecutive days when the maximum air temperature is higher than the average summer (June, July and August) maximum temperature of a baseline value +3 °C |

76    [a] Supplementary material.

## 1. INTRODUCTION

The ability of soils to sequester and store large amounts of carbon (C) is well known (e.g. Lehmann and Kleber, 2015). Soil organic carbon (SOC) stocks are crucial for maintaining soil fertility and preventing erosion and desertification, and they positively influence the provision of ecosystem services at the local as well as the global scale (e.g. Lal, 2004, 2014). For these reasons, farmers aim to establish and maintain high organic C stocks in agricultural soils, which have often been depleted trough historical land use practices (Fuchs et al., 2016; Gardi et al., 2016; Chenu et al., 2018). The continuing studies on SOC sources and biogeochemical processes in the soil environment provide key insights into climate-C feedbacks, and help prioritizing C sequestration initiatives (Gross and Harrison, 2019). In light of the climate change issue, the storage of C and additional sequestration of atmospheric C have received increasing attention recently (Rumpel et al., 2018; Whitehead et al., 2018; Lavallee et al., 2020), promoting land management, and agro-ecosystems in particular, as a key mitigation option (e.g. the '4 per mille Soils for Food Security and Climate' initiative, Minasny et al., 2017; Soussana et al., 2017). However, the slow response of SOC to changes in management and environmental factors hampers our understanding of how SOC can be increased in a sustainable manner, especially under changing climatic conditions. Long-term field experiments (LTEs), in which SOC responses have been observed over several decades, provide this information and deliver reference data on SOC content for knowledge gain and model development (Johnston and Poulton, 2018). However, LTEs are costly to maintain, and it is generally difficult to extrapolate experimental results across space and time (Debreczeni and Körschens, 2003; Mirtl et al., 2018). Simulation models play a prominent role in SOC research because they provide a mathematical framework to integrate, examine and test the understanding of SOC dynamics (Campbell and Paustian, 2015). They can also be used to extrapolate from micro- (e.g. carbohydrate production during photosynthesis) to macro-scale dynamics (e.g. global C cycling) (e.g. Gottschalk et al., 2012; Sitch et al., 2003). In particular, complex agricultural and

102 environmental models incorporate a mechanistic view of processes and system interactions, in

103 which the soil components are often represented by different, operationally defined, pools of

104 different sizes and with different properties (e.g. Parton et al., 2015). The concept of multiple C-N

105 pools represents C-N dynamics with an idealised description (Hill, 2003). The relative proportion

106 of C and N (and sometimes lignin to N ratio) in the plant residue is the primary mode to divide

107 plant inputs (from e.g. leaf litter and root exudates) into fresh litter pools, which then decompose

108 into SOC (or SOM, i.e. soil organic matter) pools, each being modelled with different residence

109 (or turnover) times, varying from months for labile products of microbial decomposition to

110 hundreds to thousands of years for organic substances with firm organic-mineral bonds (e.g. Yadav

111 and Malanson, 2007; Dungait et al., 2012). Plant material and animal manures are often modelled

112 to enter the soil environment as either readily decomposable (carbohydrate-like) or resistant (lignin

113 and cellulose-like) materials. A varying number of pools (often including inert and slow-

114 decomposing organic matter, and microbial biomass) linked by first-order equations is usually

115 simulating both C and N fluxes within and between each pool (Falloon and Smith, 2010). However,

116 different models vary considerably in the underlying assumptions and C processes in current

117 models, e.g. regarding number of pools, type of decomposition kinetics used and processes

118 regulating SOC retention (Manzoni and Porporato, 2009; Cavalli et al., 2019).

119 Each model offers a distinctive synthesis of scientific knowledge (Brilli et al., 2017) and

120 multi-model ensembles developed from several models may reduce uncertainties in biological and

121 physical outputs that occur over large scales, such as regions and continents (e.g. Rötter et al.,

122 2012; Asseng et al., 2013; Ehrhardt et al., 2018). The advantage of using ensemble estimates over

123 individual models is that caused by compensation of errors across models, and a broader

124 integration of model processes (Martre et al., 2015). It has been recommended to use model

125 ensembles for reducing uncertainties in simulations of agricultural production (Asseng et al., 2013;

126 Bassu et al., 2014; Challinor et al., 2014; Li et al., 2015; Ruane et al., 2016; Maiorano et al., 2017)

127  and other biophysical/biogeochemical outputs (Sándor et al., 2017, 2018a; Ehrhardt et al., 2018).

128  However, after the pioneering study of Smith et al. (1997), who evaluated nine SOC models using

129  12 datasets from seven LTEs, other modelling studies targeting SOC dynamics have often been

130  limited in scope. Smith et al. (2012) used four models to assess the effect on SOC of crop residues'

131  removal in 14 experiments in North America. Todd-Brown et al. (2013, 2014) performed global

132  estimates of SOC changes with 11 Earth system models. Kirschbaum et al. (2015) used one

133  simulation model and two years of eddy covariance measurements collected over an intensively

134  grazed dairy pasture in New Zealand to better understand the drivers of changes in SOC stocks.

135  Puche et al. (2019) performed a similar study in France. Using multi-model ensembles in scenario

136  studies at eight sites worldwide, Basso et al. (2018) highlighted the importance of soil feedback

137  effects (C and N) on the prediction of wheat and maize yield. We are not aware of any recent

138  model inter-comparison studies specifically assessing soil C dynamics with several models across

139  a range of experimental sites. This is a field where there is a need for standardised guidance to

140  estimate C stocks at various spatial scales (Bispo et al., 2017). A difficulty in testing and comparing

141  various models (and interpreting model outputs) lies in the interaction between soil and plant

142  processes so that any of the model-data discrepancies could be due to errors in either component

143  (e.g. Ehrmann and Ritz, 2014). A rigorous model testing and comparison would require different

144  model components, e.g. plant and soil modules, to be assessed separately. Bare-fallow plots offer

145  such an opportunity in that they are plots maintained for decades without any plant inputs. The

146  changes in SOC stocks therefore result only from decomposition processes. To assess the function

147  of soil-model components without interaction with plant processes, we conducted a model inter-

148  comparison using a dataset from long-term bare-fallow experiments where plant inputs were zero.

149  In this study, we refer to bare-fallow plots that were kept free of plants by manual and/or chemical

150  means for several decades. We used seven bare-fallow treatments included in six long-term

151  agricultural experiments (>25 years), all located in Europe (Denmark, France, Russia, Sweden and

152  United Kingdom). In these plots, the soils became progressively depleted in the more labile SOM

153  components, as they decomposed, and relatively enriched in more stable SOM (Barré et al., 2010).

154  The soil C concentrations determined at given years in these sites represented a unique opportunity

155  to follow the decay of SOC from a multi-model ensemble perspective, without any interference

156  from new plant C inputs, and conduct a multi-model ensemble comparison. The model inter-

157  comparison included 26 process-based models from an international modelling community. Some

158  models only accounted for soils  and used C input from plants as an external input where others

159  were full agro-ecosystem models that explicitly simulate plant growth and resulting C input into

160  soils. These models all simulate interactions between the soil-atmosphere continuums in different

161  ways, but for this comparison all models were run assuming no input of fresh plant-derived C,

162  allowing the comparison of just the soil components of the models.

163  Here, we assess the models, by comparing multi-decadal simulations to experimental data

164  from seven sites in Europe. The primary goal of this study was to assess the multi-model ensemble

165  in simulating SOC dynamics across bare-fallow sites in Europe. To achieve this goal, model

166  evaluation against actual measurements was performed before and after model calibration. In

167  addition, deficient areas in models and their processes were identified, paving the road for future

168  research directions.

169

170  ## 2.  MATERIALS AND METHODS

171  ### 2.1.  Simulation models

172  The ensemble of models consisted of 26 process-based models, mainly developed for crop or

173  grassland ecosystems (or focussing just on soils) and covering a broad variety of approaches (Table

174  1). While they are mostly based on first-order decay kinetics of multiple C pools (where C losses

175  are proportional to SOC stocks with additional modifiers to represent the effects of other factors),

176  ESOC1 simulates C fluxes with second-order kinetics equations based on concepts applied in

177    Schimel and Weintraub (2003) and reviewed in Wutzler and Reichstein (2008). In this case,

178    organic matter decomposition includes reactions between SOC and decomposers (i.e. a microbial

179    or enzyme pool). These different approaches depend mainly on alternative ways in which the C

180    pools are linked. For instance, MONICA is one of the most complex models, considering three

181    types of organic matter in six conceptual pools, viz. newly added organic matter, living soil

182    microbial biomass and native non-living soil organic matter, each sub-divided into fast and slowly

183    decomposing sub-pools. It simulates the turnover of C pools by applying first-order degradation

184    to each pool due to microbial growth and maintenance respiration (after Abrahamsen and Hansen,

185    2000). Then, like other models (e.g. CenW), MONICA also includes a coupled N-cycle and

186    sophisticated temperature and water-balance calculations that act as modifiers of degradation and

187    respiration rates. The decomposition rates of individual pools in such multi-pool SOC models are

188    typically controlled by vastly different reaction coefficients that can result in highly nonlinear

189    behaviour of the overall system (e.g. Caruso et al., 2018). The initial list included 34 models, but

190    eight of them were excluded from further analysis because they showed severe limitations to run

191    properly either under bare-fallow soils or under the given climate conditions. For all models,

192    estimates of SOC were compared with measured SOC data.

193 Table 1. The process-based simulation models used (model names were anonymised in the

194 reporting of simulation results using model codes from M01 to M34, from the initial list of 34

195 models, the order of models not being identical to that used in the table).

196

| Model name | Version | C pools[a] | Spin-up | URL or contact for documentation/description | References |
|---|---|---|---|---|---|
| AMG | 2 | 2 to 3 | None | https://www6.hautsdefrance.inra.fr/agroimpact/Nos-dispositifs-outils/Modeles-et-outils-d-aide-a-la-decision/AMG-et-SIMEOS-AMG/AMG-model-description | Andriulo et al. (1999); Saffih-Hdadi and Mary (2008); Clivot et al. (2019) |
| APSIM | Apsim 7.9-r4044 | 3 | None<br>Simulation from start of climate record (no additional simulation period) | http://www.apsim.info | Keating et al. (2003); Holzworth et al. (2014) |
|  | 7.10 r4158 |  | Yes |  |  |
| CANDY_CIPS | 1.0 (but always implemented in newest | 4 | None | https://www.ufz.de/export/data/2/95948_CANDY_MANUAL.pdf | Kuka, (2005); Kuka et al. (2007) |

| Model | Version | | Spin-up | Website | References |
|---|---|---|---|---|---|
| | version of CANDY 29.06.2018 | | | | |
| CCB | 2019.1.16 | 3 | None | https://www.ufz.de/index.php?en=44046 | Franko et al. (2011); Franko and Spiegel (2016); Franko and Merbach (2017) |
| Century | 4.0 | 5 to 7 | Yes | https://www2.nrel.colostate.edu/projects/century/MANUAL/html_manual/man96.html | Parton et al. (1987, 1994) |
| CenW | 4.2 | 5 | Uses an automatic spin-up routine to find equilibrium conditions under given environmental variables and specified system properties | http://www.kirschbaum.id.au/Welcome_Page.htm | Kirschbaum (1999); Kirschbaum and Paul (2002) |
| C-TOOL | 2014 | 3 | None (can be run also with spin-up) | http://envs.au.dk/fileadmin/Resources/DMU/Luft/emission/SINKS/C-TOOL_Documentation__2015_.pdf | Taghizadeh-Toosi and Olesen (2016); Taghizadeh-Toosi et al. (2014a, b, 2016) |
| Daily DayCent | 4.5 2010 / Daily DayCent 4.5 2013 | 5 to 9 | Yes | http://www.nrel.colostate.edu/projects/daycent-home.html | Parton et al. (1994, 1998); Del Grosso et al. (2001, 2002) |

Daily

DayCent

August 2014

------------------------------

4.5 2013

| DNDC | CAN | 6 | Yes<br>(10 years recommended) | http://www.dndc.sr.unh.edu | Li et al. (2012); Smith et al. (2020) |
|---|---|---|---|---|---|
| DSSAT | … | 5 | Yes,<br>20 years prior to beginning of the experiment to estimate the proportions of carbon in each organic matter pool | http://dssat.net | Jones et al. (2003); Porter et al. (2009); Gijsman et al. (2002); White et al. (2011); Thorp et al. (2012) |
| ECOSSE | 5.0.1 | 5 | None | https://www.abdn.ac.uk/staffpages/uploads/soi450/ECOSSE%20User%20manual%20310810.pdf | Smith et al. (2007, 2010a, b); Bell et al. (2010) |
| ESOC1 | 1.0 | 3 | Yes | https://doi.org/10.5281/zenodo.3539484<br>fmoyano@uni-goettingen.de | Moyano et al. (2018) |

| | | | | | |
|---|---|---|---|---|---|
| Exp | | 1 | None | - | Lorenzo Menichetti (lorenzo.menichetti@slu.se) |
| Exp + inert | | 2 | None | - | |
| ICBM | … | 2 | None | martin.bolinder@slu.se<br><br>https://www.slu.se | Andrén and Kätterer (1997); Andrén et al. (2008) |
| MONICA | 2.0.2 | 7 | None | http://monica.agrosystem-models.com | Nendel et al. (2011); Specka et al. (2016); Stella et al. (2019) |
| ORCHIDEE | 2.0 | 3 | Yes | https://vesg.ipsl.upmc.fr/thredds/fileServer/IPSLFS/orchidee/ DOXYGEN/webdoc_2425/annotated.html | Krinner et al. (2005) |
| RothC | RothC10N<br>―――<br>26.3 | 4 to 5 | None | https://www.rothamsted.ac.uk/rothamsted-carbon-model-rothc | Coleman and Jenkinson (1999); Farina et al. (2013) |
| STICS | 9.0 | 2 to 4 | None | http://www6.paca.inra.fr/stics | Brisson et al. (1998, 2003, 2008); Coucheney et al. (2015) |

| YASSO15 | 15 | 5 | Yes | https://en.ilmatieteenlaitos.fi/yasso | Tuomi et al. (2009) |

197  [a] Some models/model versions include options for varying C pools (this varying number may depend on the fact that the full

198  set of pools including fresh C can be optionally simplified in the case of bare-fallow treatments).

## 2.2. Experimental sites

We used data from a network of six long-term bare-fallow experimental sites (LTBF) in Europe (with two fields located in Askov, Denmark; Barré et al., 2010), to test the ability of the models to represent SOC dynamics. The sites were located at a range of latitudes between 48° to 59° North (Table 2; Fig. 1a), with experiments running for at least 28 years, which were used as a test bed for the models to represent SOC dynamics. Table 2 shows the main characteristics of each site and provides a brief description of the historical land use and management of the area (more details are given by Barré et al., 2010 and references therein). The documented history of the experimental sites referred to the presence of agricultural areas (grassland or cropland), without woodlands. Soil texture provides evidence of variability in soil physical properties, with a gradient of intermediate situations between the sandy loam of Askov (Denmark) and the clay loam of Ultuna (Sweden). Water relations (precipitation minus reference evapotranspiration) indicate positive climatic water balance for the two North Atlantic sites only (Askov in Denmark and Rothamsted in the United Kingdom). Mean annual temperatures vary from ~6 °C in the Sweden and Russian sites (Ultuna and Kursk, respectively) to near 11 °C in the two French sites (Grignon and Versailles). Annual air temperature amplitudes - from about 14 °C in Rothamsted to near 30 °C in Kursk - indicate that the study sites span a broad thermal gradient (Fig. 1b), which likely leads to different soil thermodynamics (e.g. Zhu et al., 2019). Two widely used metrics (aridity index and frequency of heatwaves; Sándor et al., 2017, 2018a, b) were also calculated to complete the climatic analysis of study sites (Fig. A, supplementary material).

220

221 Table 2. Long-term bare-fallow experimental sites. Table A in the supplementary material contains

222 the summary description of the experimental sites.

| | **Experimental sites (country)** | | | | | | |
|---|---|---|---|---|---|---|---|
| | **General description** | S1, S2 | S3 | S4 | S5 | S6 | S7 |
| | | Askov | Grignon | Kursk | Rothamsted | Ultuna | Versailles |
| | | (Denmark) | (France) | (Russia) | (United Kingdom) | (Sweden) | (France) |
| **Coordinates** | **Latitude** | 55.28 | 48.51 | 51.73 | 51.82 | 59.49 | 48.48 |
| | **Longitude** | 9.07 | 1.55 | 36.19 | 0.35 | 17.38 | 2.08 |
| **Soil** | **Sand/Silt/Clay (%)** | 78/12/10 (sandy loam) | 16/54/30 (silty clay loam) | 5/65/30 (silty clay loam) | 13/62/25 (silt loam) | 23/41/36 (clay loam) | 26/57/17 (silt loam) |
| | **Bulk density (Mg m$^{-3}$)** | 1.50 | 1.20 | 1.13 | 0.94 | 1.44 | 1.30 |
| | **Experimental period** *Bare-fallow years* | 1956-1985 | 1959-2007 | 1965-2001 | 1959-2008 | 1956-2007 | 1929-2008 |
| | *N. of data/replicates* | 30/4, 29/4 | 11/6 | 6/0 | 14/4 | 18/4 | 9/6 |
| | **Initial/final carbon stocks (Mg C ha$^{-1}$)** | 52.1/36.4 | 41.7/25.4 | 100.3/79.4 | 71.7/28.6 | 42.5/26.9 | 65.5/22.7 |
| **Climate[a]** | **Climate type[b]** | Dfb (humid continental) | Cfb (oceanic) | Dfb (humid continental | Cfb (oceanic) | Dfb (humid continental | Cfb (oceanic) |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | **Mean annual precipitation total (mm)** | 890 | 584 | 482 | 723 | 457 | 608 |
| | **Mean annual cumulative evaporation (mm)[c]** | 578 | 662 | 602 | 630 | 546 | 668 |
| | **Mean annual air temperature (°C)** | 7.4 | 10.7 | 6.2 | 9.4 | 6.0 | 10.7 |
| | **Mean annual air temperature range (°C)[d]** | 17.6 | 16.8 | 29.8 | 14.4 | 22.8 | 16.7 |
| **Vegetation (historical period)[e]** | **ANPP (g C m$^{-2}$ yr$^{-1}$)** | 1.7 | 1.1 | 0.9 | 1.3 | 0.9 | 1.2 |
| | **TNPP (g C m$^{-2}$ yr$^{-1}$)** | 3.3 | 2.2 | 1.7 | 2.5 | 1.7 | 2.2 |

223 [a] Climatic analysis was performed on longer periods than the experimental periods: 1956-1987/1929-2008/1944-

224 2003/1856-2006/1956-1999/1929-2008.

225 [b] Köppen-Geiger climate classification (Kottek et al., 2006).

226 [c] Mean values over the bare-fallow period. Reference evaporation was estimated based on the Thornthwaite (1948)

227 equation.

228 [d] Mean difference in temperature between the warmest and the coldest month of the year.

229 [e] Estimates of aboveground (ANPP) and total (TNPP) net primary productivity based on the precipitation levels of

230 each site, as provided by Del Grosso et al. (2008) for non-tree dominated systems.

231

232



Fig. 1. Location (a) and characterisation of the study sites (b) with respect to mean annual temperature (°C) and mean annual temperature range (°C). Details about study sites are in Table 2.

237

**2.3. Study design**

Model simulations were carried out independently by each modelling team (which included model developers and users, and field experts of soil C dynamics) on commonly formatted data using their own approaches and technical background. Harmonising calibration techniques was out of scope of the inter-comparison exercise. The SOC outputs from each model were compared to data from the study sites before and after calibration. The latter mostly focussed on parameters related to substrate use, C partitioning among pools and decomposition processes. However, rate equations for C pools often required the calibration of a large number of parameters, which are at the core of key processes responsible for differences among models in the understanding and interpretation of SOC processes (number of pools and type of decomposition kinetics used to represent C turnover). For the uncalibrated (blind test, Bln) simulations, the models were run for each site using the available data of weather and soil texture, and the initial SOC values, with no parameter adjustment other than initialisation based on historical management and land use. With this information, Bln reflects the ability of the models to simulate SOC decomposition after plant inputs has stopped, using the original parameter settings and calibration, simply by removing their components related to new C inputs. At this stage, default values were mostly used for all decomposition rates. C-pool fraction sizes were adjusted based only on C-input estimates from the information on land use prior to the establishment of the bare-fallow treatments.

After the blind simulations were completed, SOC measurements taken during the bare-fallow period were supplied to each modelling group for the calibration work. Details on management (tillage), which may have influenced the SOC dynamics before the bare-fallow treatment, were also provided to improve the initialisation process. It was requested that each modelling group adjust soil parameters to improve the simulations based on the observed data, using whatever techniques they normally use, and to document the changes. At this stage, models were split into two categories: a) with spin-up (SP) and b) without spin-up (NS). Both SP and NS

263 models require an initial estimate for SOC content and/or an adjustment of parameters towards

264 balancing the split between soil C pools. The two classes of models work in the same way using

265 information about plant residues and root growth that provide the C substrate for SOC dynamics

266 simulations. NS-type models (e.g. DNDC and RothC) use the initial measured SOC value, where

267 estimates of C inputs in the background of model runs are obtained with various methods (e.g.

268 Keel et al., 2017) in order to initialise the SOC pools, which can sometimes be calculated

269 analytically. In order to keep the legacy effect of previous land-use and past management practices,

270 in SP models (e.g. DayCent) SOC pools are routinely initialised by running the models to achieve

271 their own states of equilibrium, where change in C stocks approaches zero (e.g. Lardy et al., 2011;

272 Huntzinger et al., 2013). However, if soils are not at equilibrium (e.g. after a sudden disturbance),

273 spin-up runs may not always be valid with the risk of starting simulations with biased initial values

274 (e.g. Wutzler and Reichstein, 2007; Nemoto et al., 2017) but a fuller discussion on the "spin-up

275 problem" (Reynolds et al., 2007) is not within the scope of this paper. Carbon inputs are usually

276 estimated through sub-models calculating total net primary production (TNPP). As it was not

277 possible to derive TNPP data from local sources at each study-site, TNPP estimates were obtained

278 at each site (Table 2) based on precipitation levels according to the approach of Del Grosso et al.

279 (2008). In this way, the creation of the TNPP database used by modellers was based on an identical

280 methodology, which is widely used worldwide, though the uncertainty in quantifying productivity

281 across ecosystems is highlighted (e.g. Wieder et al., 2014).

282 The distinction between SP and NS models can appear somewhat arbitrary as virtually any

283 model with more than one C pool could be spun-up or, alternatively, a function (or analytical

284 procedures) can be used to make an initial pool partition. We refer here to common modelling

285 practice, as performed by users within the constraints imposed by packaged (operational) solutions

286 of SOC models (for which spin-up procedures may be operationally more difficult) or relying on

287 the procedure suggested by previous experience. For instance, although spin-up equilibrium runs

288  are documented for RothC (e.g. Herbst et al., 2018), it is common practice to initialise three C

289  pools for subsequent simulations through an internal routine over 10,000 years, with limited model

290  inputs including clay fraction and weather, and a pre-defined ratio of decomposable over

291  recalcitrant plant material (e.g. Xu et al., 2011; Weihermüller et al., 2013). Modellers were left to

292  choose one option or the other when both were available for use in their models (e.g. C-TOOL).

293  About 40% of the models (10 models) in the study did not use SP processes and set the initial SOC

294  values manually (using the initial SOC observation).

295      For each model category (SP and NS), two main modelling approaches were identified: site-

296  specific *versus* generic (single set of parameter values for all the sites). For the site-specific

297  approach, at each site users informed models about historical management practices and land uses

298  such as grassland or cropland (with both SP and NS models), SOC decomposition parameters (only

299  for SP models) or the partitioning of C among different soil pools (only for NS models). With the

300  generic (not site-specific) approach, model calibration was not applied separately for each

301  experimental site but simultaneously on all available multi-location datasets to find for each model

302  parameter values that would be applicable at regional scales. In this case, multi-location calibration

303  was used to capture generic model parameter values so that the models could still perform well

304  across a range of climate and management conditions in Europe (Dechow et al., 2019). Site-

305  specific and non-site-specific approaches were variously combined with factors affecting model

306  initialisation/parameterisation (Table 3) to create simulation scenarios Gen (generic), Mix (mixed)

307  and Spe (specific).

308      Scenario Mix was a regional study, where each model was calibrated simultaneously on all

309  datasets to find parameter values that would be applicable at regional scale, assuming that a single

310  calibration across sites was appropriate under the conditions explored. Fixed decomposition rate

311  parameters (but not rate modifiers) were maintained at a constant value throughout all sites (e.g.

312  the maximum passive pool decomposition rate in M25 was set to 0.003 yr$^{-1}$ at all sites), while site-

313   specific climate and soil textural conditions provided supplementary factors driving the actual

314   decomposition curve (likely in the uncalibrated blind simulations as well). In scenario Spe,

315   decomposition rates could be changed separately at each experimental site, which constrained the

316   modelling to a fitting exercise, but made it possible to explore the spatial variability of model

317   parameters. Scenario Gen ignored base histories of each site: arable crops and grasslands were not

318   distinguished, past climate conditions were disregarded, and this translated into discounting the

319   variability in the TNPP levels among sites affecting the starting SOC level.

320

321   Table 3. Modelling approaches and simulation scenarios for spin-up and no spin-up models (Gen:

322   generic; Mix: mixed; Spe: specific).

| Model category | Factors | Approaches | Calibration scenarios[a] | | |
|---|---|---|---|---|---|
| | | | Gen | Mix | Spe |
| Spin-up (SP) based models | Historical management/land use | Site-specific | | X | X |
| | | Non-site-specific | X | | |
| | Decomposition processes | Site-specific | | | X |
| | | Non-site-specific | X | X | |
| No spin-up (NS) based models | Partitioning of C pools | Site-specific | | X | X |
| | | Non-site-specific | X | | |
| | Decomposition processes | Site-specific | | | X |
| | | Non-site-specific | X | X | |

323

324   Twenty-six modelling teams participated in the blind test. At calibration stage, 17 teams

325   completed scenarios Spe and Mix, and 16 the scenario Gen. Some model packages are set to restrict

326   access to individual parameter values, which did not allow users to carry out some site-specific

327   scenarios (Mix and Spe). The same outputs were obtained with some models (e.g. RothC, DNDC),

328    which run blind and generic simulations with non-specific information like the previous land-use

329    type (arable crop or grassland) and the historical climate. When results from the blind test were

330    exactly equal to outputs from Gen scenario they were not included for further analysis. Estimated

331    and observed SOC values (Mg C ha$^{-1}$) were compared at blind test and for each calibration

332    scenario. The agreement between simulations and observations was evaluated by the inspection of

333    time series graphs and, numerically, through a set of performance metrics (Table 4) combining

334    difference- and correlation-based metrics (e.g. De Jager et al., 1994; Moriasi al., 2007;

335    Confalonieri et al., 2009; Bellocchi et al., 2002, 2010).

336

337    Table 4. Model performance metrics (P, predicted value; O, observed value; n, number of P/O

338    pairs; i, each of P/O pairs; $\overline{O}$, mean of observed values; $\overline{D}$, average of the differences between

339    predicted and observed values; $S_D$, standard deviation of the differences between estimated and

340    observed values).

| Performance metric | Equation | Unit | Value range and purpose |
|---|---|---|---|
| RRMSE, relative root mean square error (Jørgensen et al., 1986) | $$RRMSE = 100 \cdot \frac{\sqrt{\frac{\sum_{i=1}^{n}(P_i - O_i)^2}{n}}}{\overline{O}}$$ | % | 0 (optimum) to positive infinity: the closer the values are to 0, the better the model performance |
| EF, modelling efficiency (Nash and Sutcliffe, 1970) | $$EF = 1 - \frac{\sum_{i=1}^{n}(P_i - O_i)^2}{\sum_{i=1}^{n}(O_i - \overline{O})^2}$$ | - | negative infinity to 1 (optimum): the closer the values are to 1, the better the model |

| Coefficient of determination ($R^2$) of the linear regression estimates versus measurements / r, Pearson's correlation coefficient of the estimates versus measurements (Addiscott and Whitmore, 1987) | $R^2$ $$= \frac{\sum_{i=1}^{n}(P_i - O_i) \cdot (O_i - \overline{O}\_)}{\sqrt{\sum_{i=1}^{n}(P_i - \overline{P}\_)^2 \cdot \sum_{i=1}^{n}(O_i - \overline{O})^2}}$$ $$r = \sqrt{R^2}$$ | - | 0 (absence of fit of the regression line) to 1 (perfect fit of the regression line): the closer the values are to 1, the better the model<br><br>-1 (full negative correlation) to 1 (full positive correlation): the closer the values are to 1, the better the model |
| P(t), Paired Student t-test probability of means being equal | $$P(t) = Probability\left(\frac{\overline{D}}{\frac{S_D}{\sqrt{n}}}\right)$$ | - | 0 (absence of agreement) to 1 (perfect agreement): the closer the values are to 1, the better the model |
| d, index of agreement (Willmott and Wicks, 1980) | $$d = 1 - \frac{\sum_{i=1}^{n}(O_i - P_i)^2}{\sum_{i=1}^{n}(|P_i - \overline{O}| + |O_i - \overline{O}|)^2}$$ | - | 0 (absence of agreement) to 1 (perfect agreement): the closer the values are to 1, the better the model |

341

## 2.4. Multi-model and ensemble assessment

343 We first focussed on the quantification of model-data discrepancies and then assessed the

344 uncertainty of the individual models in comparison with the multi-model ensemble. The modelling

345 teams provided deterministic model simulation results according to the protocol established, which

346    meant that: 1) one run was provided for each site; 2) the spread of model results due to parameter

347    uncertainty was not specifically addressed. The latter would have dramatically increased the range

348    of model outputs used within the study and would have confounded the uncertainty in calibrated

349    parameters with the uncertainty in model structure (Wallach and Thorburn, 2017). While the

350    uncertainty in model predictions could be due to parameterisation, there is no conclusive reason

351    for regarding the model calibration from different users (something like ensemble of users within

352    ensemble of models) as being the solution to estimate uncertainty due to parameterization

353    (Confalonieri et al., 2016). Moreover, with reduction of uncertainties being mainly limited by the

354    poor quality of calibration data, rigorous calibration procedures are generally considered to be of

355    lower priority in agricultural ensemble modelling (e.g. Angulo et al., 2013; Maiorano et al., 2017).

356    As well, different calibration techniques do not seem to be primarily responsible for differences in

357    model performance (Wallach et al., 2020) and the contribution of the initialisation to the

358    uncertainty in SOC changes can be negligible compared to the uncertainty related to the model

359    itself and simulated systems characteristics (Dimassi et al., 2018). As uncertainty could not be

360    associated with any individual simulation, we focussed on the analysis of model residuals. We

361    documented the variability of the multi-model simulation exercise across two stages (blind test

362    and alternative calibration scenarios), while inspecting how the multi-model median (MMM)

363    converged to the observations. We used box-plots to compare the variability of estimates by

364    different models (with focus on multi-year averages) to the observed variability, and we

365    represented model ensembles with MMM, which has the advantage to exclude distinctly biased

366    model members with a disproportionate influence on the mean (Rodríguez et al., 2019). The

367    advantage of using MMM was established in practical studies in crop and grassland modelling but

368    also on a theoretical basis (Wallach et al., 2018).

369        We also quantified the relationship among standardised model residuals of SOC, based on

370    uncalibrated (Bln) and calibrated (Gen, Mix, Spe) simulations. Moreover, we quantified the

relationship between residuals of agro-climatic metrics (annual values): temperature amplitude, mean maximum temperature and annual precipitation. Arrays of pairwise scatterplots (scatterplot matrices) were generated with the panel plot option in the R language and environment for statistical computing ('panel.smooth', https://stat.ethz.ch/R-manual/R-devel/library/graphics/html/panel.smooth.html), which also overlaid a local non-parametric smoother curve (locally estimated scatterplot smoothing) on each plot to give some indication of trends (after Cleveland, 1979).

To explore how MMM varied with the number of models in the ensemble, we performed a calculation for each $z$-score transformed MMM, $z = \frac{MMM - \bar{O}}{sd_{obs}}$, which was obtained by dividing the multi-model data deviation from the mean of observations ($\bar{O}$) by the standard deviation of the observations ($sd_{obs}$) (Sándor et al., 2020). A $z$-score can be placed on the normal distribution curve to indicate how much it deviates from the mean of the distribution. The units of a $z$-score are $sd$ units: zero equals the mean, positive $z$-scores exceed the mean, and negative $z$-scores are less than the mean. A $z$-score allows comparisons to be made between combinations of models with different distribution characteristics, i.e. different $\bar{O}$ and $sd_{obs}$ (used here as practical descriptors of time-series central tendency and spread). As illustrated in Fig. 2, different sites occupy distinct zones in the $sd_{obs}$ versus $\bar{O}$ space. Low variability and low mean SOC observations were found at Askov (S1, S2), Grignon (S3) and Utuna (S6). The variability was higher at Rothamsted (S5) and Versailles (S7), while the mean was the highest at Kursk (S4). None of the site occupies the upper right quadrant, i.e. high variability and high mean.

Fig. 2. Standard deviation (SD) and mean of SOC observations at the study sites (details are in Table 2).

We calculated $z$-scores for all possible combinations of sets of $k$ out of $n=26$ models ($k=2, \dots n$). The minimum number of models providing plausible estimates at each site was that for which the $z$-scores lay within the ranges -1 to +1 or -2 to +2. The arbitrary choice of these thresholds was due to a conventional rule, for which values falling within 1 and 2 times the standard deviation approximate the 68% ($|z|=1$) and 95% ($|z|=2$) confidence limits of a normal distribution, respectively (after Ehrhardt et al., 2018). R software (https://cran.r-project.org) was used for statistical analysis and graphical visualization.

## 3.  RESULTS

### 3.1. Evaluation of SOC dynamics

Fig. 3 show the range of model results (represented by the shaded area) for each scenario and the multi-model median (MMM hereinafter) together with the measured values. In general, the greatest spread of model results was found under the Bln scenario, followed by the Gen scenario. In some cases, the multi-model median of Bln and Gen scenarios overestimate observations (e.g.

409    at S5, S6 and S7 sites). As expected, the tightest range of model results (simulation envelope) was

410    found with site-specific simulations. MMM simulations of Spe came closest to the observations.

411    All the MMM lines were remarkably close to the observations at sites S1, S2 and S3 (Fig. 3),

412    despite the much wider spread of the individual simulations, while the MMM at other sites differed

413    more substantially from the observations (e.g. S5, S6 and S7, Fig. 3). Overall, most of the

414    simulations (Bln, Gen and Mix) tended to overestimate the amount of SOC (e.g. S5, S6 and S7,

415    Fig. 3).

416       SOC stocks decreased under all bare-fallow sites during the investigated period. At S1, S2,

417    S3, S4 and S6 (Fig. 3) sites, the decrease in SOC stock was from minimum to moderate whereas

418    at S5 and S7 (Fig. 3) SOC loss in the top 0.20 m was more rapid, with initial SOC halved during

419    ~30 years. The decay tended to be more rapid in the first years and then the rate of loss decreased

420    (e.g. at S7 site between 1929 and 1962, Fig. 3).

421

422

Fig. 3. Temporal changes of soil organic carbon (SOC, Mg C ha$^{-1}$) observations (Observed, purple

square) and simulations: blind (Blind, blue) simulations (26 models); three calibration scenarios,

Generic (16 models, pink), Mixed (17 models, green) and Specific (17 models, grey) at all sites (

426 as in Table 2). Lines represent the multi-model median (MMM) of the simulations and shaded area

427 represents the simulation envelope.

428

429 **3.2. Ensemble performance by site**

430 Fig. 4 shows a high variability in the multi-model spread of responses at different sites. The results

431 show that Kursk (S4) soil, which stored the highest amount of SOC, 91.8 Mg C ha$^{-1}$, was

432 approximated well by the models, mainly with calibration scenario Spe, with a MMM value of

433 90.1 Mg C ha$^{-1}$. For calibration scenario Gen, some underestimation is apparent (84.2 Mg C ha$^{-1}$).

434 Site S4 had the narrowest variability in the measured values, whilst the Bln simulation and

435 calibration scenario Gen had the highest variability. Measured SOC was well estimated at S1, S2

436 and S3, including with blind simulations, despite several outlying dots, mainly with Bln and Gen

437 scenarios. The MMM tended to overestimate the measured SOC at S5 (42.5 Mg C ha$^{-1}$) and S7

438 (33.0 Mg C ha$^{-1}$) with some scenarios: Bln, S5: 56.7 Mg C ha$^{-1}$, S7: 44.49 Mg C ha$^{-1}$; Mix scenario,

439 S5: 50.0 Mg C ha$^{-1}$, S7: 35.5 Mg C ha$^{-1}$; Gen scenario, S5: 52.1 Mg C ha$^{-1}$, S7: 40.0 Mg C ha$^{-1}$.

440 On the other hand, the MMM of Gen scenarios showed the closest values to the observed median

441 at S5 and S7 (Fig. 4.).

442 Overall, with some exceptions, the MMM of calibrated runs were within the range of the

443 25$^{th}$ and 75$^{th}$ percentiles of observations. The Spe scenario provided the best MMM estimation.

444

445 **Ensemble performance by site**



447 Fig. 4. Soil organic carbon (SOC, Mg C ha$^{-1}$) at each site (as in Table 2), for blind simulations

448 (Blind, (26 models), three calibration scenarios (Mixed, 17 models; Specific and Generic, 16

449 models) and observations (Observed). For each boxplot, black horizontal lines show the multi-

450 model median. Boxes delimit the 25$^{th}$ and 75$^{th}$ percentiles. Whiskers are 10$^{th}$ and 90$^{th}$ percentiles.

451 Dots indicate outliers.

452

453 **3.3. Individual models versus multi-model ensemble**

454 The scatterplot analysis for both each model and the MMM shows that SOC estimates were

455 improved when moving from the Bln runs (Fig. 5) to the calibration Spe scenario (Fig. 6). Model

456 performances for calibration Mix and Spe scenarios also showed better simulation results than the

457 Bln simulations (see also Appendix A and Appendix B). Considering all the sites and years, the

458 predictions of some of the models (e.g. M02, M13, M22, M24 and MMM) were close to the

459 observations even for the blind level simulations (correlation coefficient >0.9, Fig. 5). Simulations

460 improved even further (correlation coefficient >0.98 for half of the models, Fig. 6) under scenario

461 Spe.

462    All the correlation coefficients of the simulations by other models also considerably improved with

463    the site-specific data and got closer to the 1:1 line. For instance, for M31, the spread of simulation

464    data in the blind simulations was mainly caused by incorrect initial SOC estimates for the different

465    sites. When the model was re-run with correctly set initial SOC amounts, the subsequent

466    drawdown of SOC over the bare-fallow period was estimated fairly well.

467    Even with blind simulations, MMM gave results in agreement with the observations ($R^2$=0.94).

468    This level of agreement was only exceeded by M22 ($R^2$=0.95) and approached by M02 ($R^2$=0.92)

469    and M13 ($R^2$=0.90). The MMM simulations continued to give the closest agreement with the

470    observations even under the full site-specific calibrations ($R^2$=0.99) with several other models

471    performing equally well (i.e. M02, M05, M09, M13, M23, M26). Overall, with some specific

472    information for model calibration, many models did remarkably well in reproducing the observed

473    patterns of SOC loss over time.

474

**Scenario Blind of SOC**

**MMM**
y = 8.56+0.859x
$R^2$ = 0.937

**M01**
y = 16+0.774x
$R^2$ = 0.729

**M03**
y = 23.1+0.593x
$R^2$ = 0.671

**M05**
y = 13.4+0.829x
$R^2$ = 0.789

**M02**
y = 3.9+0.87x
$R^2$ = 0.914

**M04**
y = 11.1+0.854x
$R^2$ = 0.833

**M06**
y = 20+0.705x
$R^2$ = 0.838

**M07**
y = -5.46+1.01x
$R^2$ = 0.757

**M09**
y = 5.25+0.917x
$R^2$ = 0.886

**M12**
y = 11.5+0.828x
$R^2$ = 0.829

**M13**
y = 5+0.924x
$R^2$ = 0.9

**M16**
y = 26.4+0.379x
$R^2$ = 0.327

**M18**
y = 11.2+0.767x
$R^2$ = 0.812

**M19**
y = 0.482+0.804x
$R^2$ = 0.687

**M20**
y = 21.6+0.732x
$R^2$ = 0.622

**M22**
y = 2.71+0.943x
$R^2$ = 0.949

**M23**
y = -6.71+1x
$R^2$ = 0.754

**M24**
y = 0.45+0.934x
$R^2$ = 0.924

**M25**
y = 5.14+0.666x
$R^2$ = 0.408

**M26**
y = 10.6+0.735x
$R^2$ = 0.855

**M27**
y = -5.26+1.17x
$R^2$ = 0.759

**M28**
y = 0.544+0.882x
$R^2$ = 0.92

**M29**
y = 27.7+0.579x
$R^2$ = 0.203

**M30**
y = 1.64+0.984x
$R^2$ = 0.817

**M31**
y = 33.8+0.808x
$R^2$ = 0.344

**M32**
y = 2.16+0.857x
$R^2$ = 0.571

**M34**
y = 39.5+-0.0717x
$R^2$ = 0.00314

SOC Simulations (Mg C ha$^{-1}$)

SOC Observations (Mg C ha$^{-1}$)

475

Fig. 5. Multi-year, multi-site comparison of individual model simulation of SOC (Mg C ha$^{-1}$): multi-model medians (MMM) from blind simulations (26 models as in Table 1) versus observations (coloured symbols represent sites as in Fig. 1).

479

Fig. 6. Multi-year, multi-site comparison of individual model simulation of SOC (Mg C ha$^{-1}$): multi-model medians (MMM) from Specific scenario simulations (17 models as in Table 1) versus observations (coloured symbols represent sites as in Fig. 1).

### 3.4. Analysis of model residuals

The plots of the discrepancy between MMM and observations (Fig. 7) as a function of time shows a limited scatter (within ±1) at each site. While Bln, Gen and Mix scenario overestimated the SOC decomposition rate at Kursk (where the highest SOC content was measured), the standardized

489   residuals were around zero at Grignon and both Askov sites during the whole of experimental

490   period. However, the departure from observations may increase over time especially with Bln and

491   Gen scenarios at some site (e.g. at Rothamsted, Ultuna, Versailles) indicating that models

492   underestimate decomposition rates after a few years/decades.

493



494

495   Fig. 7. Standardized model residuals ( $\frac{MMM-O}{sd_{obs}}$ ) over time for blind (Blind) simulations and

496   calibration scenarios Mixed, Specific and Generic at each site.

497

498   Model residuals displayed one versus the other can help establish relationships by exploring the

499   correlation of residuals from different modelling scenarios, both among them and with external

500   drivers. Residuals of blind test and calibration scenarios calculated from MMM (Fig. 8) and

501   individual models (Figs. B1-26 in the supplementary material) were correlated with the mean

502   annual climate indicators such as the precipitations, maximum temperatures and temperature

503   amplitudes. When considering the MMM, residuals of Bln were strongly correlated with Gen

504   (r=0.90) and with Mix (r=0.59) residuals, but less with Spe (r=0.25) residuals, indicating a higher

505  similarity of the first three approaches, while residuals of Spe were more correlated with those of

506  Mix (r=0.65) than of Gen (r=0.39).

507  The most prominent effect of annual climate indicators was found at the blind test stage, whose

508  residuals were negatively correlated with precipitation (r=-0.21) and positively correlated with

509  Tmax (r=0.38). Combinations of high maximum air temperature and low precipitation values may

510  thus generate greater errors in blind SOC simulations. Calibration scenario Gen did not show

511  significant correlations to climate indicators. However, calibration scenario Spe and Gen had

512  opposite correlations. The annual precipitation positively correlated with Spe residuals (r=0.26)

513  and with scenario Mix (r=0.15). Annual maximum temperature and scenario Spe negatively

514  correlated (r=-0.10). These correlations with climate indicators hint that the site-specific

515  calibration (scenario Spe) is more sensitive to precipitation than to maximum temperatures. On the

516  contrary, Bln and Gen simulation residuals showed greater sensitivity to maximum temperatures.

517  Residuals of individual models were approximately equally influenced by precipitation and

518  temperature drivers, but with differences among models and scenarios (Figs. B1-26 in the

519  supplementary material). In most of the cases, model residuals were positively correlated with

520  annual maximum temperatures and negatively correlated with annual precipitation totals (e.g.

521  M03, M09, M18, M22 for Bln). In some cases, e.g. M09 (Fig. B8 in the supplement), the

522  correlations among SOC residuals for different scenarios were both positive and negative (r values

523  ranged from -0.043 to 0.36), and even the effect of climate indicators were different (e.g. for Tmax,

524  r values ranged from -0.096 to 0.65). In other cases, e.g. M25 (Fig. B18 in the supplement), SOC

525  residuals were more similar to each other (r-values 0.17-0.80) and the effect of precipitation and

526  temperature drivers was often important (with r>0.4). It is interesting in this respect that the Spe

527  residuals had near-zero correlations with climatic drivers, showing a lesser influence of these

528  factors on model results with this scenario, whereas the Bln scenario showed some correlations

529  with Tamp (r=0.13), Tmax (r=-0.44) and precipitation (r=0.40). For M25, Gen scenario residuals

530  (Fig. B18 in the supplement) appeared unrelated with precipitation (r-value near zero), but not with

531  temperature amplitude (r=0.50) and maximum air temperature (r=-0.56).

532



533

534  Fig. 8. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of multi-model medians

535  (MMM) for blind simulations (Blind) and calibrations scenarios (Generic, Mixed and Specific as

536  in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature

537  amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother

538  curve.

539

**3.5. Minimum ensemble size**

We attempted to identify the minimum number of models required to obtain reliable results for Bln and calibration scenarios Mix, Spe and Gen (Fig. 9 and Appendix C-E). We observed that there could be large differences in the *z*-scores obtained across sites with different ensemble sizes and scenarios. Overall, Bln is characterised by greater *z*-scores than the calibration scenarios. Our analysis suggests that the ensemble size could be reduced to four models (or even fewer) at S3, S6 and S7. For the other sites (e.g. S4), only ensemble sizes of at least 9-10 models reduced *z*-scores to within the range from -2 to +2, but this number should be raised to 20 or higher to comply with the most stringent criterion of $z=|1|$. A minimum ensemble size of 9-10 models was also identified with Gen at S4 (Fig. 9), while with Mix and Spe scenarios the number of models could be reduced down to 7 and 3, respectively (up to about 14 [Gen], 8 [Mix] and 4 [Spe] to comply with $z=|1|$) (Appendix C-E).



**Generic scenarios**

552

553  Fig. 9. *z*-scores calculated with different ensemble sizes for SOC estimates obtained with Generic

554  scenario at different experimental sites. Black lines show median values. Boxes delimit the 25[th]

555  and 75[th] percentiles. Whiskers are 10[th] and 90[th] percentiles. Circles indicate outliers. Coloured

556  bands mark two critical values: *z*=|1| (light purple) and *z*=|2| (light blue).

557

558  **4.    DISCUSSION**

559  **4.1. Scenarios of ensemble SOC estimates**

560  For Bln, Mix, Gen and Spe scenarios, the overall differences between the simulated and the

561  observed initial SOC values were −0.46, +3.49, +2.40 and +1.92 Mg C ha$^{-1}$, respectively, for the

562  NS models, and +0.58, -0.29, +0.95 and -0.12 Mg C ha$^{-1}$, respectively, for the SP models. Despite

563  manually setting the initial SOC values (magnitude of first SOC observation for the simulation

564  period), the NS models mostly overestimated SOC content in the initial year of the model run. In

565  first-year estimates of the calibrated (mainly with Spe and Mix scenarios), SP models deviated less

566  from observations than NS models that overestimated SOC stocks for the first year with the

567  exception of M25 (+8.4 Mg C ha$^{-1}$ for Gen), M29 (+18.6, +21.1 and +23.7 Mg C ha$^{-1}$ for Spe, Gen

568  and Mix, respectively) and M31 (+25.2 Mg C ha$^{-1}$ for Gen). In the case of M25, the model was run

569  with a generic grassland spin-up (i.e. 7,000 years), which was applied to all sites. Thus, a generic

570  history was simulated without considering the cropping history at each site. This spin-up protocol

571  affected the simulated SOC, showing the poor ability of Gen scenario to produce results consistent

572  with observations, which questions the practicality of spin-up processes under generic calibration.

573  With M31, there was a greater difference between simulated and observed SOC values in the initial

574  simulation year and the model gave results that did not correspond to the observations at all sites

575  (Appendix F), especially under the Bln and Gen scenarios. Though M31 used the initial SOC

576  observation as default parameter, it failed to reproduce the LTBF dynamics between sites because

577  of large differences in C input to the soil from the former vegetation during the spin-up period.

578  Consequently, the starting points of the LTBF simulations differed greatly from the observations,

579  which were overestimated at S1, S2, S3 and S6, and underestimated at S4. Overall, Mix and Spe

580  calibrations showed better performance indices than the Gen scenario (Appendix F). We note,

581  however, that M13, for which the SOC pool sizes (humads and humus) were generically calibrated

582  across sites, produced low RRMSE for Gen (5.7%).

583  The improved calibration knowledge obtained with the site-specific information also improved

584  model accuracy. Moving from Bln (with knowledge of weather and soil texture, historical land use

585  and management, and initial SOC; section 2.3) to the Gen scenario, we reproduced SOC data in a

586  number of European bare-fallow experimental sites with a single set of calibrated, regional-scale

587  parameter values (regardless of the possible soil, climate and past land-use dissimilarities between

588  different sites). According to performance indicators in Appendix F, in the Bln simulations the NS

589  models performed better than the SP models. For instance, average RRMSE and EF were 19.44%

590  and 0.60, and 26.94% and 0.24, for NS and SP models, respectively. Compared to the Bln scenario,

591  the discrepancy between the measured and estimated SOC values under the Gen scenario was

592  slightly reduced with NS models and increased with SP models. Multi-site calibration can be

593  characterised by lower uncertainty than site-specific calibration, because more data contribute to

594  the calibration process (e.g. Minunno et al., 2014; Ma et al., 2015). The availability of a variety of

595  detailed data from multiple sites thus offers the possibility of a genuine multi-location calibration

596  of the model, assuming that a single calibration across sites is appropriate. The limit of the Gen

597  scenario calibration was that it did not make it possible to explore the spatial variability of model

598  parameters. The latter was done with scenarios Mix and Spe, for which a basic requisite is that

599  model parameters are not hard-coded but configuration files are left open to the users. From Gen

600  to Mix, parameters describing initial values of each pool were determined separately for each site.

601  Moving from Mix to Spe, the decomposition parameters became site-specific. Hence, modellers

602  needed to invest increasingly more knowledge (and more time-demanding calibration effort) than

603    in Gen. Under these conditions, the improvement of simulations in SP models was evident (up to

604    70% for some indicators, e.g. RRMSE and EF). On the contrary, NS models only had a slight

605    improvement in accuracy of simulations from Bln (RRMSE=21.5%; EF=0.58) to Mix

606    (RRMSE=18.6%, EF=0.55) or Gen (RRMSE=20.5%; EF=0.45). In our analysis, the two types of

607    models (NS and SP) appear to be suitable for different sets of data. NS-type models, in most cases,

608    can perform well even when data are limited to climate, initial C and historic land use, while SP

609    models generally benefit from the availability of more detailed data. All metrics related to the

610    performance of the SP models were improved with calibration. There were some differences in

611    model performance among the sites, but site-specific soil or climatic conditions cannot easily

612    explain such differences.

613    Overall, across the seven LTEs and using simulated and observed SOC data at the end of the

614    experimental period we observe that the greatest and least differences from observations were

615    approximately +14.3% with Bln and +2.2% with Spe (Fig. 10). The Gen scenario achieved almost

616    half the error (+8.9%) of is closest competitor, i.e. the Bln scenario. More than one-third of the

617    Bln-scenario error is achievable with the Mix scenario (+4.0%).

618



619

Fig. 10. Multi-site averages (vertical bars) and standard deviations (vertical lines) of observed and estimated (ensemble multi-model median) values of SOC (Mg C ha$^{-1}$) in the last year of the experimental period. The ensemble modelling was applied with blind simulations (Blind) and calibration scenarios (Mixed, Specific and Generic as in Table 3).

This study has shown that it is difficult to define an *a priori* criterion that could be used to select a subset of models that would perform better than others would. In terms of the minimum number of models required to obtain reliable results, our study indicates that the suggested minimum ensemble size (~10 models) proposed by Martre et al. (2015) for crop growth could be a reference also when model ensembles are implemented to blindly simulate SOC in bare-fallow soils, which can be reduced down to 3-4 models with a site-specific calibration. These sizes are lower than that found by Sándor et al. (2020) to provide reliable C-flux estimates in croplands and grasslands (i.e. ~13 models). While the current study applied the same methodology as Sándor et al. (2020), but as the present study focuses on one output variable only, SOC, evaluated in simplified systems (bare-fallow soils), its relative ease of simulation offers great advantages for scenario analyses in the absence of vegetation cover and plant residues, nor farming practices (only occasional tillage operations occurred at some sites and were considered by models which can simulate this option). This is reflected in the several *z*-scores within the range of -2 and +2, as obtained with a limited number of models, showing that reduced ensemble sizes can satisfactorily estimate the SOC content in bare-fallow systems, mainly when site-specific calibration is possible. However, our analysis of the Russian site (S4), which had low observed variability and high mean ($sd_{obs}$=6.9, $\bar{O}$=91.8 Mg C ha$^{-1}$), is challenging because it showed that model ensembles that are too small might not always guarantee sufficient accuracy in SOC estimates of C-rich soils. An application to the peatlands located on the Mid-Russian Upland (e.g. Shumilovskikh et al., 2018) should thus be considered with caution.

645

**4.2. Possibilities for model inaccuracies**

647 We presented an approach that uses a correlation matrix (with graphical representation) to account

648 for possible correlations between Bln, Mix, Gen and Spe residuals and, additionally, climatic

649 factors (mean air temperature amplitude, maximum air temperature and precipitation total). This

650 residual analysis helps find correlations among alternative scenarios, which might indicate

651 comparable scenarios in which error propagation within models is similar, though the way of error

652 propagation cannot be easily retrieved from the correlation matrix. This is the case of Bln, Gen

653 and Mix, whose residuals are highly correlated, while the weak correlations between Spe and other

654 scenarios highlight the distinct behaviour of the latter. This analysis can also help find correlations

655 between the SOC output and external drivers, and thus suggest additional predictors that may need

656 to be included in the models (e.g. Medlyn et al., 2005). This need emerged especially when specific

657 models were run under Bln, Gen and Mix scenarios, for which some correlations (r>|0.4|) were

658 obtained between model residuals and drivers of thermal and moisture conditions. A weaker but

659 significant correlation (r=0.26, p=0.02) was also obtained between Spe residuals and precipitation.

660 These correlations indicate some limitations related to the response functions of SOC

661 decomposition to soil temperature and soil moisture, though the relative uncertainties of our model

662 ensemble are attenuated by the presence in the models of physical and chemical processes that

663 explain the intra- and inter-annual variability of SOC. We add that such biophysical conditions

664 affect the microbial activity (e.g. Blagodatskaya and Kuzyakov, 2008; Guenet et al., 2010; Wutzler

665 and Reichstein, 2013), and care should be taken when extrapolating our results over long time

666 frames (especially without locally calibrated models, Fig. 7) if no corroborating field evidence for

667 long-term decay rates can be obtained (e.g. on how models are dealing such situations in which

668 microbes become increasingly C limited as no new C input by plants occurs; Kuhry and Vitt,

669 1996).

670

## 5.        CONCLUSIONS AND FUTURE DIRECTIONS

This paper on SOC modelling offers a tentative answer to the questions about: (i) whether and to what extent an ensemble of models performs better than single models, (ii) the minimum ensemble size that is required to reduce the error below a given threshold, and (iii) the set of data required to prepare and substantiate ensemble estimates. This study presents a framework for interpretation of model performance and uncertainties obtained with a set of process-based biogeochemical models (individually and in an ensemble) simulating soil C contents in bare-fallow experimental systems at a variety of European sites. One of the features of SOC modelling today is the huge amount and variety of models available. Although our analysis was not taking into account all sources of uncertainty (e.g. the influence of the unique choices made by modellers), it enabled the integration of several modelling teams into an ensemble protocol. Classifying and comparing different approaches have revealed great model diversity, and is the basis for the development of dedicated ensemble protocols. In this model inter-comparison, the need to accommodate challenges experienced by modellers (including C pools of different nature, and optional initialisation and calibration procedures) was reflected in the co-creation (with modellers and data providers) of alternative calibration scenarios (Mix, Gen, Spe). As far as we are aware, no previous multi-model inter-comparison studies have examined differences in such calibration scenarios or differences between models with or without spin-up.

In our study, we did not aim to identify the best model(s) for simulating SOC dynamics for bare-fallows and no probability of success was assigned to prove the suitability of using one model rather than another. Overall, we showed that a calibration scenario with generic system knowledge was adequate for providing sufficiently reliable output, but additional site-specific knowledge can further improve results under certain circumstances. This is operationally relevant because the effort required to gather calibration data might no longer be feasible for modelling scenarios

695     moving from single sites to increasingly larger spatial scales. Site-specific calibration could help

696     refine model estimates. However, geographical locations have characteristics (e.g. soil and climate

697     conditions, past history) that require specific model structures and local optimisation, and the

698     application of models may be limited by the ability to provide representative parameter values.

699     Soil-C model inter-comparisons including more models and experimental data from other regions

700     should be continued to improve our ability to simulate biogeochemical processes with acceptable

701     accuracy. Additional assessments are also recommended to complete the analysis of model

702     behaviour in the long term (like thousands of years) with constant inputs. While the various models

703     evaluated here did not include all available modelling approaches used to simulate soil C

704     dynamics, the present model inter-comparison was large compared to other studies. As such, it is

705     a distinct improvement over previously published quantitative approaches because it represents a

706     reasonable sub-population of common and current approaches. In this, we offer a method to allow

707     a broad ensemble of models to be implemented using existing datasets and current modelling

708     practices. Overall, this multi-model ensemble sets a precedent for key progress in soil C modelling

709     because it provides essential information about SOC modelling and opens a path to a more in-

710     depth analysis of the response of individual models and their uncertainties against soil and climate

711     drivers. Now that we have examined SOC decomposition in-depth without the difficulties of C

712     input uncertainties, a similar modelling study should be conducted on LTEs that examine both

713     plant derived C inputs as well as C inputs from manures and other organic materials recycled in

714     agroecosystems. How simulation models compare under such conditions is important for

715     improving our ability to evaluate and achieve climate C goals. With increasing availability of data

716     and computational resources, there are many opportunities for the SOC modelling community to

717     enrich its offering and to keep up with evolving methodologies, which would significantly increase

718     transparency of the underpinning science and modelling practice. A number of recent actions are

719     ongoing under the guidance of international initiatives such as the European Joint Programme

720   (EJP) on Soil (https://projects.au.dk/ejpsoil). Started in 2020, the EJP-Soil is undertaking a detailed

721   inventory of models and all available data sources (e.g. world soil maps, satellite images,

722   downscaled weather data), and appears as an ideal arena to facilitate the exchange of information

723   and to further explore SOC model developments and practice.

724   **ACKNOWLEDGEMENTS**

744

**AUTHOR CONTRIBUTIONS**

R. Farina, R. Sándor and G. Bellocchi coordinated the study, contributed to its design, conducted the analysis of data and produced the first draft of the manuscript. P. Smith, C. Chenu, F. Ehrhardt, M. A. Bolinder, C. Nendel and J.-F. Soussana contributed to the design of the study and the writing of the manuscript. M. Abdalla, J. Álvaro-Fuentes, M. A. Bolinder, L. Brilli, H. Clivot, M. De Antoni, C. Di Bene, C. D. Dorich, F. Ferchaud, N. Fitton, R. Francaviglia, U. Franko, D. Giltrap, B. B. Grant, B. Guenet, M. T. Harrison, M. U. F. Kirschbaum, K. Kuka, L. Kulmala, J. Liski, M. J. McGrath, E. Meier, L. Menichetti, F. Moyano, N, Reibold, A. Shepherd, W. N. Smith, T. Stella, A. Taghizadeh-Toosi and E. Tsutskikh performed the model calibrations and runs.

C. Dorich, L. Bechini, L. Menichetti, R. Francaviglia, S. Recous, W. Smith, F. Ferchaud, H. Clivot, M. A. Bolinder, W. Smith, A. Taghizadeh-Toosi, L. Brilli, R. Farina, G. Bellocchi, T. Stella and U. Franko discussed and decided upon the modelling scenarios at the CN-MIP final meeting (Rome, 6-7 June 2018). C. Dorich prepared a detailed protocol for second-stage simulations. Those interested in the details of the modelling process are encouraged to contact authors.

**REFERENCES**

Abrahamsen, P., & Hansen, S. (2000). Daisy: an open soil-crop-atmosphere system model. *Environmental Modelling & Software*, **15**, 313-330. https://doi.org/10.1016/S1364-8152(00)00003-7

Addiscott, T. M., & Whitmore, A. P. (1987). Computer simulation of changes in soil mineral nitrogen and crop nitrogen during autumn, winter and spring. *Journal of Agricultural Science*, **109**, 141-157. https://doi.org/10.1017/S0021859600081089

Andrén, O., & Kätterer, T. (1997). ICBM: The introductory carbon balance model for exploration of soil carbon balances. *Ecological Applications*, **7**, 1226-1236. https://doi.org/10.1890/1051-0761(1997)007[1226:ITICBM]2.0.CO;2

770  Andrén, O., Kätterer, T., Karlsson, T., & Eriksson, J. (2008). Soil C balances in Swedish

771    agricultural soils 1990-2004, with preliminary projections. *Nutrient Cycling in Agroecosystems*,

772    **81**, 129–144. https://doi.org/10.1007/s10705-008-9177-z

773  Andriulo, A., Mary, B., & Guerif, J. (1999). Modelling soil carbon dynamics with various cropping

774    sequences    on    the    rolling    pampas.    *Agronomie*,    **19**,    365–377.

775    https://doi.org/10.1051/agro:19990504

776  Angulo, C., Rötter, R., Lock, R., Enders, A., Fronzek, S., & Ewert, F. (2013). Implication of crop

777    model calibration strategies for assessing regional impacts of climate change in Europe.

778    *Agricultural    and    Forest    Meteorology*,    **170**,    32–46.

779    https://doi.org/10.1016/j.agrformet.2012.11.017

780  Asseng, S., Ewert, F., Rosenzweig, C., Jones, J. W., Hatfield, J. L., Ruane, A., … Wolf, J. (2013).

781    Uncertainty in simulating wheat yields under climate change. *Nature Climate Change*, **3**, 827–

782    832. https://doi.org/10.1038/nclimate1916

783  Barré, P., Eglin, T., Christensen, B. T., Ciais, P., Houot, S., Kätterer, T., ... Chenu, C. (2010).

784    Quantifying and isolating stable soil organic carbon using long-term bare fallow experiments.

785    *Biogeosciences*, **7**, 3839-3850. https://doi.org/10.5194/bg-7-3839-2010

786  Basso, B., Dumont, B., Maestrini, B., Shcherbak, I., Robertson, G. P., Porter, J. R., … Rosenzweig,

787    C. (2018). Soil organic carbon and nitrogen feedbacks on crop yields under climate change.

788    *Agricultural and Environmental Letters*, **3**, 180026. https://doi.org/10.2134/ael2018.05.0026

789  Bassu, S., Brisson, N., Durand, J. L., Boote, K., Lizaso, J., Jones, J. W., … Waha, K., 2014. How

790    do various maize crop models vary in their responses to climate change factors? *Global Change*

791    *Biology*, **20**, 2301–2320. https://doi.org/10.1111/gcb.12520

792  Bellocchi, G., Acutis, M., Fila, G., & Donatelli, M. (2002). An indicator of solar radiation model

793    performance based on a fuzzy expert system. *Agronomy Journal*, **94**, 1222-1233.

794    https://doi.org/10.2134/agronj2002.1222

795  Bellocchi, G., Rivington, M., Donatelli, M., & Acutis, M. (2010). Validation of biophysical

796      models: issues and methodologies. A review. *Agronomy for Sustainable Development*, **30**, 109-

797      130. https://doi.org/10.1051/agro/2009001

798  Bispo, A., Andersen, L., Angers, D. A., Bernoux, M., Brossard, M., Cécillon, L., … Eglin, T.K.

799      (2017). Accounting for carbon stocks in soils and measuring GHGs emission fluxes from soils:

800      do we have the necessary standards? Frontiers in Environmental Science, **12 July 2017**.

801      https://doi.org/10.3389/fenvs.2017.00041

802  Blagodatskaya, E., & Kuzyakov, Y. (2008). Mechanisms of real and apparent priming effects and

803      their dependence on soil microbial biomass and community structure: critical review. *Biology*

804      *and Fertility of Soils*, **45**, 115–131. https://doi.org/10.1007/s00374-008-0334-y

805  Brilli, L., Bechini, L., Bindi, M., Carozzi, M., Cavalli, D., Conant, R., … Bellocchi, G. (2017).

806      Review and analysis of strengths and weaknesses of agro-ecosystem models for simulating C

807      and   N   fluxes.   *Science   of   the   Total   Environment*,   **598**,   445-470.

808      https://doi.org/10.1016/j.scitotenv.2017.03.208

809  Brisson, N., Mary, B., Ripoche, D., Jeuffroy, M. H., Ruget, F., Nicollaud, B., … Delécolle, R.

810      (1998). STICS: a generic model for the simulation of crops and their water and nitrogen

811      balances. I. Theory and parameterization applied to wheat and corn. *Agronomie*, **18**, 311–346.

812      https://doi.org/10.1051/agro:19980501

813  Brisson, N., Gary, C., Justes, E., Roche, R., Mary, B., Ripoche, D., … Sinoquet, H. (2003). An

814      overview of the crop model STICS. *European Journal of Agronomy*, **18**, 309-332.

815      https://doi.org/10.1016/S1161-0301(02)00110-7

816  Brisson, N., Launay, M., Mary, B., & Baudoin, N. (2008). Conceptual basis, formalizations and

817      parameterization of the STICS crop model. Paris (France): Editions Quae.

818  Campbell, E. E., & Paustian, K. (2015). Current developments in soil organic matter modeling and

819     the expansion of model applications: a review. *Environmental Research Letters*, **10**, 123004.

820     https://doi.org/10.1088/1748-9326/10/12/123004

821  Caruso, T., De Vries, F., Bardgett, R. D., & Lehmann, J. (2018). Soil organic carbon dynamics

822     matching ecological equilibrium theory. *Ecology and Evolution*, **8**, 11169-11178.

823     https://doi.org/10.1002/ece3.4586

824  Cavalli, D., Bellocchi, G., Corti, M., Gallina, P. M., & Bechini, L. (2019). Sensitivity analysis of

825     C and N modules in biogeochemical crop and grassland models following manure addition to

826     soil. *European Journal of Soil Science*, **70**, 833-846. https://doi.org/10.1111/ejss.12793

827  Challinor, A., Martre, P., Asseng, S., Thornton, P., & Ewert, F. (2014). Making the most of climate

828     impacts ensembles. *Nature Climate Change*, **4**, 77-80. https://doi.org/10.1038/nclimate2117

829  Chenu, C., Angers, D. A., Barré, P., Derrien, D., Arrouays, D., & Balesdent, J. (2018). Increasing

830     organic stocks in agricultural soils: Knowledge gaps and potential innovations. *Soil and Tillage*

831     *Research*, **188**, 41-52. https://doi.org/10.1016/j.still.2018.04.011

832  Cleveland, W.S. (1979). Robust locally weighted regression and smoothing scatterplots. J. Am.

833     Stat. Assoc. 74, 829-836. https://doi.org/10.1080/01621459.1979.10481038

834  Clivot, H., Mouny, J. C., Duparque, A., Dinh, J. L., Denoroy, P., Houot, S., … Mary, B. (2019).

835     Modeling soil organic carbon evolution in long-term arable experiments with AMG model.

836     *Environmental       Modelling       &       Software*,       **118**,       99-113.

837     https://doi.org/10.1016/j.envsoft.2019.04.004

838  Coleman, K., & Jenkinson, D.S. (1999). RothC-26.3 - A model for the turnover of carbon in soil:

839     model description and Windows user guide. Harpenden (UK): Lawes Agricultural Trust.

840  Confalonieri, R., Acutis, M., Bellocchi, G., & Donatelli, M. (2009). Multi-metric evaluation of the

841     models WARM, CropSyst, and WOFOST for rice. *Ecological Modelling*, **220**, 1395-1410.

842     https://doi.org/10.1016/j.ecolmodel.2009.02.017

843    Confalonieri, R., Orlando, F., Paleari, L., Stella, T., Gilardelli, C., Movedi, E., ... Acutis, M.

844        (2016). Uncertainty in crop model predictions: what is the role of users? *Environmental*

845        *Modelling & Software*, **81**, 165-173. https://doi.org/10.1016/j.envsoft.2016.04.009

846    Coucheney, E., Buis, S., Launay, M., Constantin, J., Mary, B., García de Cortázar-Atauri, I., …

847        Léonard, J. (2015). Accuracy, robustness and behavior of the STICS soil–crop model for plant,

848        water and nitrogen outputs: Evaluation over a wide range of agro-environmental conditions in

849        France. *Environmental Modelling & Software*, **64**, 177-190.

850        https://doi.org/10.1016/j.envsoft.2014.11.024

851    De Jager, J.M. (1994). Accuracy of vegetation evaporation ratio formulae for estimating final

852        wheat yield. *Water SA*, **20**, 307-314. Retrieved from

853        https://journals.co.za/content/waters/20/4/AJA03784738_2194

854    Debreczeni, K., & Körschens, M. (2003). Long-term field experiments of the world. *Archives of*

855        *Agronomy and Soil Science*, **49**, 465-483. https://doi.org/10.1080/03650340310001594754

856    Dechow, R., Franko, U., Kätterer, T., & Kolbe, H. (2019). Evaluation of the RothC model as a

857        prognostic tool for the prediction of SOC trends in response to management practices on arable

858        land. *Geoderma*, **337**, 463-478. https://doi.org/10.1016/j.geoderma.2018.10.001

859    Del Grosso, S. J., Parton, W. J., Mosier, A. R., Hartman, M. D., Brenner, J., Ojima, D. S., &

860        Schimel, D. S. (2001). Simulated interaction of carbon dynamics and nitrogen trace gas fluxes

861        using the DayCent model. In M. J. Shaffer, L. Ma, & S. Hansen (Eds.), *Modeling carbon and*

862        *nitrogen dynamics for soil management* (pp. 303-332). Boca Raton: CRC Press.

863    Del Grosso, S., Ojima, D., Parton, W., Mosier, A., Peterson, G., & Schimel, D. (2002). Simulated

864        effects of dryland cropping intensification on soil organic matter and greenhouse gas exchanges

865        using the DAYCENT ecosystem model. *Environmental Pollution*, **1**, S75-S83.

866        https://doi.org/10.1016/S0269-7491(01)00260-3

867     Del Grosso, S., Parton, W., Stohlgren, T., Zheng, D., Bachelet, D., Prince, S., … Olson, R. (2008).

868     Global potential net primary production predicted from vegetation class, precipitation, and

869     temperature. *Ecology*, **89**, 2117-2126. https://doi.org/10.1890/07-0850.1

870     Dimassi, B., Guenet, B., Saby, N. P. A., Munoz, F., Bardy, M., Millet, F., & Martin, M. P. (2018).

871     The impacts of CENTURY model initialization scenarios on soil organic carbon dynamics

872     simulation in French long-term experiments. *Geoderma*, **311**, 25-36.

873     https://doi.org/10.1016/j.geoderma.2017.09.038

874     Dungait, J. A. J., Hopkins, D. W., Gregory, A. S., & Whitmore, A. P. (2012). Soil organic matter

875     turnover is governed by accessibility not recalcitrance. *Global Change Biology*, **18**, 1781-1796.

876     https://doi.org/10.1111/j.1365-2486.2012.02665.x

877     Ehrhardt, F., Soussana, J.-F., Bellocchi, G., Grace, P., Mcauliffe, R., Recous, S., … Zhang, Q.

878     (2018). Assessing uncertainties in crop and pasture ensemble model simulations of productivity

879     and $N_2O$ emissions. *Global Change Biology*, **24**, e603-e616. https://doi.org/10.1111/gcb.13965

880     Ehrmann, J., & Ritz, K. (2014). Plant: soil interactions in temperate multi-cropping production

881     systems. *Plant and Soil*, **376**, 1-29. https://doi.org/10.1007/s11104-013-1921-8

882     Falloon, P., & Smith, P. (2010). Modelling soil carbon dynamics. In W. L. Kutsch, M. Bahn, & A.

883     Heinemeyer (Eds.), *Soil carbon dynamics: An integrated methodology* (pp. 221-244).

884     Cambridge: Cambridge University Press.

885     Farina, R., Coleman, K., & Whitmore, A. P. (2013). Modification of the RothC model for

886     simulations of soil organic C dynamics in dryland regions. *Geoderma*, **200-201**, 18-30.

887     https://doi.org/10.1016/j.geoderma.2013.01.021

888     Franko, U., Kolbe, H., Thiel, E., & Liess, E. (2011). Multi-site validation of a soil organic matter

889     model for arable fields based on generally available input data. *Geoderma*, **166**, 119-134.

890     https://doi.org/10.1016/j.geoderma.2011.07.019

891    Franko, U., & Spiegel, H. (2016). Modeling soil organic carbon dynamics in an Austrian long-

892        term tillage field experiment. *Soil and Tillage Research*, **156**, 83-90.

893    Franko, U., & Merbach, I. (2017). Modelling soil organic matter dynamics on a bare fallow

894        Chernozem soil in Central Germany. *Geoderma*, **303**, 93-98.

895        https://doi.org/10.1016/j.geoderma.2017.05.013

896    Fuchs, R., Schulp, C. J. E., Hengeveld, G. M., Verburg, P. H., Clevers, J. G. P. W., Schelhaas, M.-

897        J., & Herold, M. (2016). Assessing the influence of historic net and gross land changes on the

898        carbon fluxes of Europe. *Global Change Biology*, **22**, 2526-2539.

899        https://doi.org/10.1111/gcb.13191

900    Gardi, C., Visioli, G., Conti, F. D., Scotti, M., Menta, C., & Bodini, A. (2016). High Nature Value

901        Farmland: assessment of soil organic carbon in Europe. Frontiers in Environmental Science, 21

902        June 2016. https://doi.org/10.3389/fenvs.2016.00047

903    Gijsman, A. J., Hoogenboom, G., Parton, W. J., & Kerridge, P. C. (2002). Modifying DSSAT crop

904        models for low-input agricultural systems using a soil organic matter-residue module from

905        CENTURY. *Agronomy Journal*, **94**, 462-474. https://doi.org/10.2134/agronj2002.4620

906    Gottschalk, P., Smith, J. U., Wattenbach, M., Bellarby, J., Stehfest, E., Arnell, N., … Smith, P.

907        (2012). How will organic carbon stocks in mineral soils evolve under future climate? Global

908        projections using RothC for a range of climate change scenarios. *Biogeosciences*, **9**, 3151-3171.

909        https://doi.org/10.3390/soilsystems3020028

910    Gross C. D., & Harrison, R. B. (2019). The case for digging deeper: soil organic carbon storage,

911        dynamics, and controls in our changing world. *Soil Systems*, **3**, 28.

912        https://doi.org/10.3390/soilsystems3020028

913    Guenet, B., Neill, C., Bardoux, G., & Abbadie, L. (2010). Is there a linear relationship between

914        priming effect intensity and the amount of organic matter input? *Applied Soil Ecology*, **46**, 436–

915        442. https://doi.org/10.1016/j.apsoil.2010.09.006

916 Herbst, M., Welp, G., Macdonald, A., Jate, M., Hädicke, A., Scherer, H., … Vanderborght, J.
917 (2018). Correspondence of measured soil carbon fractions and RothC pools for equilibrium and
918 non-equilibrium states. *Geoderma*, **314**, 37-46.
919 https://doi.org/10.1016/j.geoderma.2017.10.047

920 Hill, M. J. (2003). Generating generic response signals for scenario calculation of management
921 effects on carbon sequestration in agriculture: approximation of main effects using CENTURY.
922 *Environmental Modelling & Software*, **18**, 899-913. https://doi.org/10.1016/S1364-
923 8152(03)00054-9

924 Holzworth, D. P., Huth, N. I., deVoil, P. G., Zurcher, E. J., Herrmann, N. I., McLean, G., …
925 Keating, B. A. (2014). APSIM - Evolution towards a new generation of agricultural systems
926 simulation. *Environmental Modelling & Software*, **62**, 327-350.
927 https://doi.org/10.1016/j.envsoft.2014.07.009

928 Huntzinger, D. N., Schwalm, C., Michalak, A. M., Schaefer, K., King, A. W., Wei, Y., … Zhu, Q.
929 (2013). The North American Carbon Program Multi-scale synthesis and Terrestrial Model
930 Intercomparison Project-Part 1: Overview and experimental design. *Geoscientific Model*
931 *Development*, **6**, 2121-2133. https://doi.org/10.5194/gmd-6-2121-2013

932 Johnston, A. E., & Poulton, P. R. (2018). The importance of long-term experiments in agriculture:
933 their management to ensure continued crop production and soil fertility; the Rothamsted
934 experience. *European Journal of Soil Science*, **69**, 113-125. https://doi.org/10.1111/ejss.12521

935 Jones, J. W., Hoogenboom, G., Porter, C. H., Boote, K. J., Batchelor, W. D., Hunt, L. A., …
936 Ritchie, J. T. (2003). The DSSAT cropping system model. *European Journal of Agronomy*, **18**,
937 235–265. https://doi.org/10.1016/S1161-0301(02)00107-7

938 Jørgensen, S. E., Kamp-Nielsen, L., Christensen, T., Windolf-Nielsen, J., & Westergaard, B.
939 (1986). Validation of a prognosis based upon a eutrophication model. Ecological Modelling,
940 **35**, 165-182. https://doi.org/10.1016/0304-3800(86)90024-4

941    Keating, B. A., Carberry, P. S., Hammer, G. L., Probert, M. L., Robertson, M. J., Holzworth, D.,

942    … Smith, C. J. (2003). An overview of APSIM, a model designed for farming systems

943    simulation. *European Journal of Agronomy*, **18**, 267-288. https://doi.org/10.1016/S1161-

944    0301(02)00108-9

945    Keel, S. G., Leifeld, J., Mayer, J., Taghizadeh-Toosi, A., and Olesen, J. E. (2017). Large

946    uncertainty in soil carbon modelling related to method of calculation of plant carbon input in

947    agricultural systems. *European Journal of Soil Science*, **68**, 953-863.

948    https://doi.org/10.1111/ejss.12454

949    Kirschbaum, M.U.F. (1999). CenW, a forest growth model with linked carbon, energy, nutrient

950    and water cycles. *Ecological Modelling*, **118**, 17–59. https://doi.org/10.1016/S0304-

951    3800(99)00020-4

952    Kirschbaum, M. U. F., Rutledge, S., Kuijper, I. A., Mudge, P. L., Puche, N., Wall, A. M., …

953    Campbell, D. I. (2015). Modelling carbon and water exchange of a grazed pasture in New

954    Zealand constrained by eddy covariance measurements. *Science of the Total Environment*, **512-

955    513**, 273-286. https://doi.org/10.1016/j.scitotenv.2015.01.045

956    Kirschbaum, M. U. F., & Paul, K. I. (2002). Modelling carbon and nitrogen dynamics in forest

957    soils with a modified version of the CENTURY model. *Soil Biology & Biochemistry*, **34**, 341-

958    354. https://doi.org/10.1016/S0038-0717(01)00189-4

959    Kottek, M., Grieser, J., Beck, C., Rudolf, B., & Rubel, F. (2006). World map of the Köppen-Geiger

960    climate classification updated. *Meteorologische Zeitschrift*, **15**, 259-263.

961    https://doi.org/10.1127/0941-2948/2006/0130

962    Krinner, G., Viovy, N., de Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P., … Colin

963    Prentice, I. (2005). A dynamic global vegetation model for studies of the coupled atmosphere-

964    biosphere system. *Global Biogeochemical Cycles*, **19**, GB1015.

965    https://doi.org/10.1029/2003GB002199

966  Kuhry, P., & Vitt, D.H. (1996). Fossil carbon/nitrogen ratios as a measure of peat decomposition.

967  *Ecology*, **77**, 271–275. https://doi.org/10.2307/2265676

968  Kuka, K. (2005). Modellierung des Kohlenstoffhaushaltes in Ackerböden auf der Grundlage

969  bodenstrukturabhängiger Umsatzprozesse. PhD thesis, Martin-Luther-University Halle-

970  Wittenberg.                                     Retrieved                                     from

971  https://gepris.dfg.de/gepris/projekt/5247578?context=projekt&task=showDetail&id=5247578

972  & (in German)

973  Kuka, K., Franko, U., & Rühlmann, J. (2007) Modelling the impact of pore space distribution on

974  carbon           turnover.           *Ecological           Modelling*,           **208**,           295–306.

975  https://doi.org/10.1016/j.ecolmodel.2007.06.002

976  Lal, R. (2004). Soil carbon sequestration impacts on global climate change and food security.

977  *Science*, **304**, 1623-1626. https://doi.org/10.1126/science.1097396

978  Lal, R. (2014). Soil conservation and ecosystem services. *International Soil and Water*

979  *Conservation Research*, **2**, 36-47. https://doi.org/10.1016/S2095-6339(15)30021-6

980  Lardy, R., Bellocchi, G., & Soussana, J.-F. (2011). A new method to determine soil organic carbon

981  equilibrium.         *Environmental         Modelling         &         Software*,         **26**,         1759-1763.

982  https://doi.org/10.1016/j.envsoft.2011.05.016

983  Lavallee, J. M., Soong, J. L., & Cotrufo, M. F. (2020). Conceptualizing soil organic matter into

984  particulate and mineral-associated forms to address global change in the 21$^{st}$ century. *Global*

985  *Change Biology*, **26**, 261-273. https://doi.org/10.1111/gcb.14859

986  Lehmann, J., & Kleber, M. (2015). The contentious nature of soil organic matter. *Nature*, **528**, 60-

987  68. https://doi.org/10.1038/nature16069

988  Li, C., Salas, W., Zhang, R., Krauter, C., Rotz, A., & Mitloehner, F. (2012). Manure-DNDC: a

989  biogeochemical process model for quantifying greenhouse gas and ammonia emissions from

990      livestock manure systems. *Nutrient Cycling in Agroecosystems*, **93**, 163-200.

991      https://doi.org/10.1007/s10705-012-9507-z

992      Li, T., Hasegawa, T., Yin, X., Zhu, Y., Boote, K., Adam, M., … Bouman, B. (2015). Uncertainties

993      in predicting rice yield by current crop models under a wide range of climatic conditions. *Global*

994      *Change Biology*, **21**, 1328-1341. https://doi.org/10.1111/gcb.12758

995      Ma, S., Lardy, R., Graux, A.-I., Ben Touhami, H., Klumpp, K., Martin, R., Bellocchi, G. (2015).

996      Regional-scale analysis of carbon and water cycles on managed grassland systems.

997      *Environmental Modelling & Software*, **72**, 356-371.

998      https://doi.org/10.1016/j.envsoft.2015.03.007

999      Maiorano, A., Martre, P., Asseng, S., Ewert, F., Müller, C., Rötter, R. P., … Zhu, Y. (2017). Crop

1000      model improvement reduces the uncertainty of the response to temperature of multi-model

1001      ensembles. *Field Crops Research*, **202**, 5-20. https://doi.org/10.1016/j.fcr.2016.05.001

1002      Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: Theory and models

1003      across scales. *Soil Biology & Biochemistry*, **41**, 1355-1379.

1004      https://doi.org/10.1016/j.soilbio.2009.02.031

1005      Martre, P., Wallach, D., Asseng, S., Ewert, F., Jones, J.W., Rotter, R.P., … Wolf, J. (2015).

1006      Multimodel ensembles of wheat growth: Many models are better than one. *Global Change*

1007      *Biology*, **21**, 911-925. https://doi.org/10.1111/gcb.12768

1008      Medlyn, B. E., Robinson, A. P., Clement, R., & McMurtrie, R. E. (2005). On the validation of

1009      models of forest $CO_2$ exchange using eddy covariance data: some perils and pitfalls. *Tree*

1010      *Physiology*, **25**, 839-857. https://doi.org/10.1093/treephys/25.7.839

1011      Minasny, B., Malone, B. P., McBratney, A. B., Angers, D. A., Arrouays, D., Chambers, A., …

1012      Winowiecki, L. (2017). Soil carbon 4 per mille. *Geoderma*, **292**, 59–86.

1013      https://doi.org/10.1016/j.geoderma.2017.01.002

1014  Minunno, F., Peltoniemi, M., Launiainen, S., & Mäkelä, A. (2014). Integrating ecosystems

1015  measurements from multiple eddy-covariance sites to a simple model of ecosystem process -

1016  are there possibilities for a uniform model calibration? *Geophysical Research Abstracts*, **16**,

1017  EGU2014-10706-3.                              Retrieved                              from

1018  https://meetingorganizer.copernicus.org/EGU2014/orals/14065

1019  Mirtl, M., Borer, E. T., Djukic, I., Forsius, M., Haubold, H., Hugo, W., Jourdane, J., … Haase, P.

1020  (2018). Genesis, goals and achievements of long-term ecological research at the global scale: a

1021  critical review of ILTER and future directions. *Science of the Total Environment*, **626**, 1439-

1022  1462. https://doi.org/10.1016/j.scitotenv.2017.12.001

1023  Moriasi, D., Arnold, J., Van Liew, M., Bingner, R., Harmel, R., & Veith, T. (2007). Model

1024  evaluation guidelines for systematic quantification of accuracy in watershed simulations.

1025  *Transactions of the ASABE*, **50**, 885-900. https://doi.org/10.13031/2013.23153

1026  Moyano, F. E., Vasilyeva, N., & Menichetti, L. (2018). Diffusion limitations and Michaelis–

1027  Menten kinetics as drivers of combined temperature and moisture effects on carbon fluxes of

1028  mineral soils. *Biogeosciences*, **15**, 5031–5045. https://doi.org/10.5194/bg-15-5031-2018

1029  Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I - a

1030  discussion of principles. *Journal of Hydrology*, **10**, 282-290. https://doi.org/10.1016/0022-

1031  1694(70)90255-6

1032  Nemoto, R., Klumpp, K., Coleman, K., Dondini, M., Goulding, K., Hastings, A., … Smith, P.

1033  (2016). Soil organic carbon (SOC) equilibrium and model initialisation methods: an application

1034  to the Rothamsted Carbon (RothC) model. *Environmental Modeling & Assessment*, **22**, 215-

1035  229.

1036  Nendel, C., Berg, M., Kersebaum, K. C., Mirschel, W., Specka, X., Wegehenkel, M., … Wieland,

1037  R. (2011). The MONICA model: Testing predictability for crop growth, soil moisture and

1038    nitrogen    dynamics.    *Ecological    Modelling*,    **222**,    1614–1625.

1039    https://doi.org/10.1016/j.ecolmodel.2011.02.018

1040    Parton, W. J., Del Grosso, S., Plante, A. F., Adair, E. C., & Lutz, S. M. (2015). Modeling the

1041    dynamics of soil organic matter and nutrient cycling. In E. A. Paul (Ed.), *Soil microbiology,*

1042    *ecology and biochemistry, 4th edition* (pp. 505-537). Amsterdam: Elsevier Academic Press.

1043    Parton, W. J., Hartman, M., Ojima, D., & Schimel, D. (1998). DAYCENT and its land surface

1044    submodel: description and testing. *Global and Planetary Change*, **19**, 35-48.

1045    https://doi.org/10.1016/S0921-8181(98)00040-X

1046    Parton, W. J., Schimel, D. S., & Cole, C.V., & Ojima, D. S. (1987). Analysis of factors controlling

1047    soil organic matter levels in Great Plains grasslands. Soil Science Society of America Journal,

1048    **51**, 1173–1179. https://doi.org/10.2136/sssaj1987.03615995005100050015x

1049    Parton, W. J., Schimel, D. S., Ojima, D. S., & Cole, C. V. (1994). A general model for soil organic

1050    matter dynamics: sensitivity to litter chemistry, texture and management. In R. B. Bryant & R.

1051    W. Arnold (Eds.), *Quantitative modeling of soil forming processes* (pp. 147–167). Madison,

1052    WI (USA): SSSA Spec. Pub. 39. ASA, CSSA and SSSA.

1053    Porter, C. H., Jones, J. W., Adiku, S., Gijsman, A. J., Gargiulo, O., & Naab, J. B. (2009). Modeling

1054    organic carbon and carbon-mediated soil processes in DSSAT v4.5. *Operational Research*, **10**,

1055    247-278. https://doi.org/10.1007/s12351-009-0059-1

1056    Puche, N. J. B., Senapati, N., Flechard, C. R., Klumpp, K., Kirschbaum, M. U. F, & Chabbi, A.

1057    (2019). Modelling carbon and water fluxes of managed grasslands: comparing flux variability

1058    and net carbon budgets between grazed and mowed systems. *Agronomy*, **9**, 183.

1059    https://doi.org/10.3390/agronomy9040183

1060    Reynolds, K. M., Thomson, A. J., Köhl, M., Shannon, M. A., Ray, D., & Rennolls, K. (2007).

1061    Sustainable forestry: from monitoring and modelling to knowledge management and policy

1062    science. Wallingford: CAB International.

1063    Rodríguez, A., Ruiz-Ramos, M., Palosuo, T., Carter, T. R., Fronzek, S., Lorite, I. J., … Rötter, R.

1064    P. (2019). Implications of crop model ensemble size and composition for estimates of

1065    adaptation effects and agreement of recommendations. *Agricultural and Forest Meteorology*,

1066    **15**, 351-362. https://doi.org/10.1016/j.agrformet.2018.09.018

1067    Rötter, R. P., Palosuo, T., Kersebaum, K. C., Angulo, C., Bindi, M., Ewert, F., … Trnka, M.

1068    (2012). Simulation of spring barley yield in different climatic zones of Northern and Central

1069    Europe – A comparison of nine crop models. *Field Crops Research*, **133**, 23–36.

1070    https://doi.org/10.1016/j.fcr.2012.03.016

1071    Ruane, A. C., Hudson, N. I., Asseng, S., Camarrano, D., Ewert, F., Martre, P., … Wolf, J. (2016).

1072    Multi-wheat-model ensemble responses to interannual climate variability. *Environmental*

1073    *Modelling & Software*, **81**, 86-101. https://doi.org/10.1016/j.envsoft.2016.03.008

1074    Rumpel, C., Amiraslani, F., Koutika, L. S., Smith, P., Whitehead, D., & Wollenberg, E. (2018).

1075    Put more carbon in soils to meet Paris climate pledges. *Nature*, 564, 32-34.

1076    https://doi.org/10.1038/d41586-018-07587-4

1077    Saffih-Hdadi, K., & Mary, B. (2008). Modeling consequences of straw residues export on soil

1078    organic    carbon.    *Soil    Biology    &    Biochemistry*,    **40**,    594–607.

1079    https://doi.org/10.1016/j.soilbio.2007.08.022

1080    Sándor, R., Barcza, Z., Acutis, M., Doro, L., Hidy, D., Köchy, M., … Bellocchi, G. (2017). Multi-

1081    model simulation of soil temperature, soil water content and biomass in Euro-Mediterranean

1082    grasslands: Uncertainties and ensemble performance. *European Journal of Agronomy*, **88**, 22-

1083    40. https://doi.org/10.1016/j.eja.2016.06.006

1084    Sándor, R., Ehrhardt, F., Brilli, L., Carozzi, M., Recous, S., Smith, P., … Bellocchi, G. (2018a).

1085    The use of biogeochemical models to evaluate mitigation of greenhouse gas emissions from

1086    managed    grasslands.    *Science    of    the    Total    Environment*,    **642**,    292-306.

1087    https://doi.org/10.1016/j.scitotenv.2018.06.020

1088    Sándor, R., Ehrhardt, F., Grace, P., Recous, S., Smith, P., Snow, V., … Bellocchi, G. (2020).

1089    Ensemble modelling of carbon fluxes in grasslands and croplands. *Field Crops Research*, **252**,

1090    107791. https://doi.org/10.1016/j.fcr.2020.107791

1091    Sándor, R., Picon-Cochard, C., Martin, R., Louault, F., Klumpp, K., Borras, D., & Bellocchi, G.,

1092    (2018b). Plant acclimation to temperature: Developments in the Pasture Simulation model.

1093    *Field Crops Research*, **222**, 238-255. https://doi.org/10.1016/j.fcr.2017.05.030

1094    Schimel, J. P., & Weintraub, M. N. (2003). The implications of exoenzyme activity on microbial

1095    carbon and nitrogen limitation in soil: a theoretical model. *Soil Biology & Biochemistry*, **35**,

1096    549–563. https://doi.org/10.1016/S0038-0717(03)00015-4

1097    Shumilovskikh, L. S., Novenko, E., & Giesecke, T. (2018). Long-term dynamics of the East

1098    European forest-steppe ecotone. *Journal of Vegetation Science*, **29**, 416-426.

1099    https://doi.org/10.1111/jvs.12585

1100    Sitch, S., Smith, B., Prentice, I. C., Arneth, A., Bondeau, A., Cramer, W., … Venevsky, S. (2003).

1101    Evaluation of ecosystem dynamics, plant geography and terrestrial carbon cycling in the LPJ

1102    dynamic global vegetation model. *Global Change Biology*, **9**, 161-185.

1103    https://doi.org/10.1046/j.1365-2486.2003.00569.x

1104    Smith, J., Gottshcalk, P., Bellarby, J., Chapman, S., Lilly, A., Towers, W., … Smith, P. (2010a).

1105    Estimating changes in national soil carbon stocks using ECOSSE – a new model that includes

1106    upland organic soils. Part I. Model description and uncertainty in national scale simulations of

1107    Scotland. *Climate Research*, **45**, 179-192. https://doi.org/10.3354/cr00899

1108    Smith, J., Gottschalk, P., Bellarby, J., Chapman, S., Lilly, A., Towers, W., … Smith, P. (2010b).

1109    Estimating changes in national soil carbon stocks using ECOSSE - a new model that includes

1110    upland organic soils. Part II. Application in Scotland. *Climate Research*, **45**, 193-205.

1111    https://doi.org/10.3354/cr00902

1112    Smith, P., Smith, J., Flynn, H., Killham, K., Rangel-Castro, I., Foereid, B., … Falloon, P., 2007.

1113    ECOSSE: Estimating Carbon in Organic Soils - Sequestration and Emissions. Final Report.

1114    SEERAD Report, 166 pp. Retrieved from http://nora.nerc.ac.uk/id/eprint/2233

1115    Smith, P., Smith, J. U., Powlson, D. S., McGill, W. B., Arah, R. M., Chertov, O. G., … Whitmore,

1116    A. P. (1997). A comparison of the performance of nine soil organic matter models using datasets

1117    from seven long-term experiments. *Geoderma*, **81**, 153-225. https://doi.org/10.1016/S0016-

1118    7061(97)00087-6

1119    Smith, W. N., Grant, B. B., Campbell, C. A., McConkey, B. G., Desjardins, R. L., Kröbel, R. &

1120    Malhi, S. S. (2012). Crop residue removal effects on soil carbon: Measured and inter-model

1121    comparisons.    *Agriculture,    Ecosystems    &    Environment*,    **161**,    27-38.

1122    https://doi.org/10.1016/j.agee.2012.07.024

1123    Smith, W. N., Grant, B., Qi, Z., He, W., VanderZaag, A., Drury, C. F., & Helmers, M. (2020).

1124    Development of the DNDC model to improve soil hydrology and incorporate mechanistic tile

1125    drainage: A comparative analysis with RZWQM2. *Environmental Modelling & Software*, **123**,

1126    104577. https://doi.org/10.1016/j.envsoft.2019.104577

1127    Soussana, J.-F., Lutfalla, S., Ehrhardt, F., Rosenstock, T. S., Lamanna, C., Havlik, P., … Lal, R.

1128    (2017). Matching policy and science: Rationale for the '4 per 1000 - soils for food security and

1129    climate'    initiative.    *Soil    and    Tillage    Research*,    **188**,    3-15.

1130    https://doi.org/10.1016/j.still.2017.12.002

1131    Specka, X., Nendel, C., Hagemann, U., Pohl, M., Hoffmann, M., Barkusky, D., … van Oost, K.

1132    (2016). Reproducing $CO_2$ exchange rates o a crop rotation at contrasting terrain positions using

1133    two    different    modelling    approaches.    *Soil    and    Tillage    Research*,    **156**,    219–229.

1134    https://doi.org/10.1016/j.still.2015.05.007

1135  Stella, T., Mouratiadou, I., Gaiser, T., Berg-Mohnicke, M., Wallor, E., Ewert, F., & Nendel, C.

1136  (2019). Estimating the contribution of crop residues to soil organic carbon conservation.

1137  Environmental Research Letters 14, 094008. https://doi.org/10.1088/1748-9326/ab395c

1138  Taghizadeh–Toosi, A., Christensen, B. T., Hutchings, N. J., Vejlin, J., Kätterer, T., Glendining,

1139  M., & Olesen, J. E. (2014a). C-TOOL: A simple model for simulating whole-profile carbon

1140  storage in temperate agricultural soils. *Ecological Modelling*, **292**, 11-25.

1141  https://doi.org/10.1016/j.ecolmodel.2014.08.016

1142  Taghizadeh-Toosi, A., Olesen, J. E., Kristensen, K., Elsgaard, L., Østergaard, H. S., Lægdsmand,

1143  M., … Christensen, B. T. (2014b). Changes in carbon stocks of Danish agricultural mineral

1144  soils between 1986 and 2009. *European Journal of Soil Science*, **65**, 730-740.

1145  https://doi.org/10.1111/ejss.12169

1146  Taghizadeh-Toosi, A., & Olesen, J. E. (2016). Modelling soil organic carbon in Danish agricultural

1147  soils suggests low potential for future carbon sequestration. *Agricultural Systems*, **145**, 83-89.

1148  https://doi.org/10.1016/j.agsy.2016.03.004

1149  Taghizadeh-Toosi, A., Christensen, B. T., Glendining, M., & Olesen, J. E. (2016). Consolidating

1150  soil carbon turnover models by improved estimates of belowground carbon input. *Scientific*

1151  *Reports*, **6**, 32568. https://doi.org/10.1038/srep32568

1152  Thornthwaite, C. W. (1948). An approach toward a rational classification of climate. *Geographical*

1153  *Review*, **38**, 55-94. https://doi.org/10.2307/210739

1154  Thorp, K. R., White, J. W., Porter, C. H., Hoogenboom, G., Nearing, G. S., & French, A. N. (2012).

1155  Methodology to evaluate the performance of simulation models for alternative compiler and

1156  operating system configurations. *Computers and Electronics in Agriculture*, **81**, 62-71.

1157  https://doi.org/10.1016/j.compag.2011.11.008

1158  Todd-Brown, K. E. O., Randerson, J. T., Post, W. M., Hoffman, F. M., Tarnocai, C., Schuur, E.

1159  A. G., & Allison, S. D. (2013). Causes of variation in soil carbon simulations from CMIP5

1160    Earth system models and comparison with observations. *Biogeosciences*, **10**, 1717–1736.

1161    https://doi.org/10.5194/bg-10-1717-2013

1162 Todd-Brown, K. E. O., Randerson, J. T., Hopkins, F., Arora, V., Hajima, T., Jones, C., … Allison,

1163    S. D. (2014). Changes in soil organic carbon storage predicted by Earth system models during

1164    the 21st century. *Biogeosciences*, **11**, 2341–2356. https://doi.org/10.5194/bg-11-2341-2014

1165 Tuomi, M., Thum, T., Järvinen, H., Fronzek, S., Berg, B., Harmon, M., … Liski, J. (2009). Leaf

1166    litter decomposition - Estimates of global variability based on Yasso07 model. *Ecological*

1167    *Modelling*, **220**, 3362-3371. https://doi.org/10.1016/j.ecolmodel.2009.05.016

1168 Wallach, D., Martre, P., Liu, B., Asseng, S., Ewert, F., Thonburn, P.J., … Zhang, Z. (2018). Multi-

1169    model ensembles improve predictions of crop-environment-management interactions. *Global*

1170    *Change Biology*, **24**, 5072-5083. https://doi.org/10.1111/gcb.14411

1171 Wallach, D., Palosuo, T., Thorburn, P., Seidel, S. J., Gourdain, E., Asseng, S., … Zhu, Y. (2020).

1172    How well do crop models predict phenology, with emphasis on the effect of calibration?

1173    *bioRxiv*, March 30, 2020. https://doi.org/10.1101/708578

1174 Wallach, D., & Thorburn, P. J. (2017). Estimating uncertainty in crop model predictions: Current

1175    situation and future prospects. *European Journal of Agronomy*, **88**, A1-A7.

1176    https://doi.org/10.1016/j.eja.2017.06.001

1177 Weihermüller, L., Graf, A., Herbst, M., & Vereecken, H. (2013). Simple pedotransfer functions to

1178    initialize reactive carbon pools of the RothC model. *European Journal of Soil Science*, **64**, 567-

1179    575. https://doi.org/10.1111/ejss.12036

1180 White, J. W., Hoogenboom, G., Kimball, B. A., & Wall, G. W. (2011). Methodologies for

1181    simulating impacts of climate change on crop production. *Field Crops Research*, **124**, 357-368.

1182    https://doi.org/10.1016/j.fcr.2011.07.001

1183 Whitehead, D., Schipper, L. A., Pronger, J., Moinet, G. Y., Mudge, P. L., Pereira, R. C., … Camps-

1184    Arbestain, M. (2018). Management practices to reduce losses or increase soil carbon stocks in

1185     temperate grazed grasslands: New Zealand as a case study. *Agriculture, Ecosystems &*

1186     *Environment*, **265**, 432-443. https://doi.org/10.1016/j.agee.2018.06.022

1187    Wieder, W. R., Boehnert, J., & Bonan, G. B. (2014). Evaluating soil biogeochemistry

1188     parameterizations in Earth system models with observations. *Global Biogeochemical Cycles*,

1189     **28**, 211-222. https://doi.org/10.1002/2013GB004665

1190    Willmott, C. J., & Wicks, D. E. (1980). An empirical method for the spatial interpolation of

1191     monthly precipitation within California. *Physical Geography*, **1**, 59-73.

1192     https://doi.org/10.1080/02723646.1980.10642189

1193    Wutzler, T., & Reichstein, M. (2007). Soils apart from equilibrium - consequences for soil carbon

1194     balance modelling. *Biogeosciences*, **4**, 125-136. https://doi.org/10.5194/bg-4-125-2007

1195    Wutzler, T., & Reichstein, M. (2008). Colimitation of decomposition by substrate and

1196     decomposers - a comparison of model formulations. *Biogeosciences*, **5**, 749–759.

1197     https://doi.org/10.5194/bg-5-749-2008

1198    Wutzler, T., & Reichstein, M. (2013). Priming and substrate quality interactions in soil organic

1199     matter models. *Biogeosciences*, **10**, 2089–2103. https://doi.org/10.5194/bg-10-2089-2013

1200    Xu, X., Wen L., & Kiely, G. (2011). Modeling the change in soil organic carbon of grassland in

1201     response to climate change: Effects of measured versus modelled carbon pools for initializing

1202     the Rothamsted Carbon model. *Agriculture, Ecosystems & Environment*, **140**, 372-381.

1203     https://doi.org/10.1016/j.agee.2010.12.018

1204    Yadav, V., & Malanson, G. (2007). Progress in soil organic matter research: litter decomposition,

1205     modelling, monitoring and sequestration. *Progress in Physical Geography*, **31**, 131-154.

1206     https://doi.org/10.1177/0309133307076478Zhu, D., Ciais, P., Krinner, G., Maignan, F., Puig,

1207     A.J., & Hugelius, G. (2019). Controls of soil organic matter on soil thermal dynamics in the

1208     northern high latitudes. *Nature Communications*, **10**, 3172. https://doi.org/10.1038/s41467-

1209     019-11103-1

1210 **Appendix A**

1211 Multi-year, multi-site comparison of individual model simulation of SOC (Mg C ha$^{-1}$): multi-

1212 model medians (MMM) from Mix scenario simulations (17 models) versus observations.

1213 (coloured symbols represent sites as in Fig. 1).

1214



1215

**Appendix B**

Multi-year, multi-site comparison of individual model simulation of SOC (Mg C ha$^{-1}$): multi-

model medians (MMM) from Gen scenario simulations (16 models) versus observations.

(coloured symbols represent sites as in Fig. 1).

1220



Scenario Generic of SOC

1221

1222

**Appendix C**

1224   *z*-scores calculated with different ensemble sizes for SOC estimates obtained with Bln scenario at

1225   different experimental sites. Black lines show median values. Boxes delimit the 25[th] and 75[th]

1226   percentiles. Whiskers are 10[th] and 90[th] percentiles. Circles indicate outliers. Coloured bands mark

1227   two critical values: *z*=|1| (light purple) and *z*=|2| (light blue).



1228

1229

1230

1231

1232    **Appendix D**

1233    *z*-scores calculated with different ensemble sizes for SOC estimates obtained with Mix scenario at

1234    different experimental sites. Black lines show median values. Boxes delimit the 25$^{th}$ and 75$^{th}$

1235    percentiles. Whiskers are 10$^{th}$ and 90$^{th}$ percentiles. Circles indicate outliers. Coloured bands mark

1236    two critical values: *z*=|1| (light purple) and *z*=|2| (light blue).



**Mixed scenarios**

1237

1238

**Appendix E**

*z*-scores calculated with different ensemble sizes for SOC estimates obtained with Spe scenario at

different experimental sites. Black lines show median values. Boxes delimit the 25[th] and 75[th]

percentiles. Whiskers are 10[th] and 90[th] percentiles. Circles indicate outliers. Coloured bands mark

two critical values: *z*=|1| (light purple) and *z*=|2| (light blue).



1244

1245

1246

**Appendix F**

1248 Individual and multi-model ensemble (MMM) performance metrics (as in Table 4) for blind (Bln)

1249 and calibration scenarios (Mix, Spe and Gen as in Table 3) across sites. Red (italic) and blue (bold)

1250 numbers indicate the worst and best performances by metric, respectively.

| Performance metric | Scenario | M01 | M02 | M03 | M04 | M05 | M06 | M07 | M09 | M12 | M13 | M16 | M18 | M19 | M20 | M22 | M23 | M24 | M25 | M26 | M27 | M28 | M29 | M30 | M31 | M32 | M34 | MMM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $R^2$ | Bln | 0.73 | 0.92 | 0.67 | 0.83 | 0.79 | 0.86 | 0.76 | 0.89 | 0.83 | 0.90 | 0.33 | 0.81 | 0.69 | 0.63 | **0.95** | 0.76 | 0.92 | 0.41 | 0.86 | 0.76 | 0.92 | *0.21* | 0.82 | 0.35 | 0.57 | 0.80 | 0.94 |
| | Gen | NA | 0.39 | NA | NA | NA | NA | NA | 0.79 | NA | **0.97** | 0.90 | NA | NA | 0.56 | 0.87 | NA | 0.89 | 0.09 | 0.86 | 0.93 | 0.91 | 0.82 | 0.93 | *~0.00* | NA | 0.85 | 0.95 |
| | Mix | NA | 0.91 | NA | 0.90 | 0.91 | NA | NA | 0.89 | NA | **0.99** | NA | NA | 0.83 | *0.41* | 0.98 | 0.56 | 0.94 | 0.49 | **0.99** | 0.95 | NA | 0.91 | 0.84 | 0.87 | NA | 0.82 | 0.97 |
| | Spe | 0.97 | **0.99** | NA | 0.98 | **0.99** | NA | NA | **0.99** | NA | **0.99** | 0.96 | NA | NA | 0.96 | 0.98 | **0.99** | NA | 0.91 | **0.99** | 0.97 | *0.88* | 0.93 | 0.98 | 0.94 | NA | NA | **0.99** |

d

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bln | 0.88 | 0.97 | 0.84 | 0.93 | 0.90 | 0.89 | 0.90 | 0.97 | 0.93 | 0.97 | 0.71 | 0.94 | 0.85 | 0.79 | **0.99** | 0.89 | 0.97 | 0.73 | 0.95 | 0.91 | 0.95 | 0.59 | 0.95 | *0.52* | 0.85 | 0.93 | 0.98 |
| Gen | NA | 0.71 | NA | NA | NA | NA | NA | 0.93 | NA | **0.99** | 0.97 | NA | NA | 0.66 | 0.96 | NA | 0.97 | 0.53 | 0.95 | 0.97 | 0.97 | 0.81 | 0.97 | *0.23* | NA | 0.94 | 0.98 |
| Mix | NA | 0.97 | NA | 0.96 | 0.97 | NA | NA | 0.97 | NA | **~1.00** | NA | NA | 0.89 | *0.69* | **~1.00** | 0.79 | 0.98 | 0.81 | **~1.00** | 0.98 | NA | 0.76 | 0.96 | 0.96 | NA | 0.93 | 0.99 |
| Spe | 0.99 | **~1.00** | NA | **~1.00** | **~1.00** | NA | NA | **~1.00** | NA | **~1.00** | 0.99 | NA | NA | 0.99 | **~1.00** | 0.99 | NA | 0.97 | **~1.00** | 0.99 | 0.95 | *0.76* | 0.99 | 0.98 | NA | NA | **~1.00** |

RRMSE (%)

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bln | 24.1 | 10.9 | 28.0 | 18.6 | 21.9 | 21.9 | 23.1 | 12.5 | 17.7 | 11.8 | 28.6 | 15.5 | 27.2 | 33.1 | **7.9** | 25.4 | 11.0 | 36.6 | 14.0 | 24.0 | 14.4 | 48.4 | 16.3 | *69.1* | 27.7 | 16.3 | 10.4 |
| Gen | NA | 30.8 | NA | NA | NA | NA | NA | 17.9 | NA | **5.7** | 11.5 | NA | NA | 51.3 | 14.0 | NA | 12.1 | 49.4 | 14.5 | 12.7 | 10.9 | 37.9 | 12.4 | *92.1* | NA | 15.8 | 10.6 |
| Mix | NA | 11.0 | NA | 12.6 | 11.5 | NA | NA | 11.7 | NA | **3.8** | NA | NA | 23.3 | 45.6 | 4.4 | 29.0 | 8.9 | 33.0 | 4.2 | 9.4 | NA | *46.5* | 14.4 | 13.4 | NA | 15.9 | 7.2 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Spe | 6.5 | 3.4 | NA | 5.0 | **3.2** | NA | NA | 3.8 | NA | 3.8 | 8.2 | NA | NA | 6.7 | 4.4 | 5.0 | NA | 14.5 | 4.1 | 6.2 | 14.9 | *46.2* | 5.5 | 8.7 | NA | NA | **3.2** |

P(t)

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bln | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | *~0.00* | **0.64** | 0.02 | *~0.00* | *~0.00* | 0.31 | *~0.00* | *~0.00* | *~0.00* | 0.45 | 0.05 | *~0.00* | *~0.00* | 0.13 | *~0.00* | *~0.00* | 0.01 | *~0.00* |
| Gen | NA | *~0.00* | NA | NA | NA | NA | NA | *~0.00* | NA | 0.13 | **0.17** | NA | NA | *~0.00* | *~0.00* | NA | 0.08 | 0.04 | *~0.00* | *~0.00* | 0.06 | *~0.00* | *~0.00* | *~0.00* | NA | *~0.00* | *~0.00* |
| Mix | NA | *~0.00* | NA | *~0.00* | *~0.00* | NA | NA | 0.55 | NA | 0.31 | NA | NA | *~0.00* | *~0.00* | **0.76** | *~0.00* | 0.54 | *~0.00* | 0.31 | *~0.00* | NA | *~0.00* | 0.24 | *~0.00* | NA | *~0.00* | 0.49 |
| Spe | 0.46 | 0.99 | NA | 0.06 | 0.03 | NA | NA | 0.85 | NA | 0.34 | *~0.00* | NA | NA | 0.12 | 0.93 | *~0.00* | NA | *~0.00* | **~1.00** | 0.29 | *~0.00* | *~0.00* | *~0.00* | 0.68 | NA | NA | 0.83 |

EF

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bln | 0.52 | 0.90 | 0.49 | 0.72 | 0.60 | 0.60 | 0.56 | 0.87 | 0.74 | 0.88 | 0.33 | 0.80 | 0.39 | 0.09 | **0.95** | 0.47 | 0.90 | -0.11 | 0.84 | 0.53 | 0.83 | -0.93 | 0.78 | *-2.95* | 0.37 | 0.78 | 0.91 |
| Gen | NA | 0.22 | NA | NA | NA | NA | NA | 0.73 | NA | 0.97 | 0.89 | NA | NA | -1.17 | 0.84 | NA | 0.88 | -0.49 | 0.83 | 0.87 | 0.90 | -0.19 | 0.87 | *-6.00* | NA | 0.79 | **0.93** |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mix | NA | 0.90 | NA | 0.87 | 0.89 | NA | NA | 0.89 | NA | **0.99** | NA | NA | 0.55 | -0.72 | 0.98 | 0.31 | 0.93 | 0.34 | **0.99** | 0.93 | NA | *-0.78* | 0.83 | 0.85 | NA | 0.79 | 0.97 |
| Spe | 0.97 | **0.99** | NA | 0.98 | **0.99** | NA | NA | **0.99** | NA | **0.99** | 0.94 | NA | NA | 0.96 | 0.98 | 0.98 | NA | 0.87 | **0.99** | 0.97 | 0.82 | *-0.76* | 0.97 | 0.94 | NA | NA | **0.99** |

1251

Supplementary material of

**"Ensemble modelling, uncertainty and robust predictions of organic carbon in long-term bare-fallow soils"**

Table A. Description of the long-term bare-fallow experimental sites.

| Experimental sites (country) | | | | | |
|---|---|---|---|---|---|
| S1, S2 Askov (Denmark) | S3 Grignon (France) | S4 Kursk (Russia) | S5 Rothamsted (United Kingdom) | S6 Ultuna (Sweden) | S7 Versailles (France) |
| It has been under cultivation since around 1800. The site was open to mixed heath and grasslands with scattered deciduous shrubs. It was historically used for occasional haymaking and light grazing. The field was frequently tilled to 0.2 m depth was fertilized with 70 kg N ha$^{-1}$ yr$^{-1}$ before 1973 and with 100 kg N ha$^{-1}$ yr$^{-1}$ thereafter. The experiment was adjacent to the B3- and B4-fields (blocks). SOC measurements were obtained from four replicates over 29 years of observations. Weeds were hand-removed. | The field was unmanaged grassland since before 1875. The field was tilled twice a year, during autumn (November) and spring (March) to 0.25 m depth. Weeds were removed by hand weeding and herbicide treatments. SOC measurements were obtained from six replicates over 48 years of observation. | Originally a steppe, the field was brought into cultivation around 200 years before 1964. Before the start of the fallow period, the field was under five-year rotation (maize-alfalfa-potato-winter wheat-pea), tilled twice a year to 0.22 to 0.25 m depth. The fallow period was maintained by mechanical and chemical weed removal. SOC measurements were obtained from one replicate and 36 years of observations. | The field was a managed grassland since 1838. Under that treatment, the field was tilled two to four times per year to 0.20 to 0.22 m depth. The fallow period was maintained by occasional weed removal through herbicides. SOC measurements were obtained from four replicates over 49 years of observations. | The field was used for crop cultivation for several centuries. Then, five to 10 years prior to LTBF conversion, the field was under nitrogen fertilisation and straw return to soil. The field was manually tilled with a spade once a year in autumn to 0.2 m depth. During the LBTF treatment, weeds were mostly mechanically removed in spring and throughout the growing season, and only once (in 1991) chemically removed. Measurements were obtained from four replicates over 51 years of observations. | The LTBF experiment forms part of a large experiment (42 plots). The field was a mixture of unmanaged grassland and forest before the 17$^{th}$ century, followed by unmanaged grassland. No information was available on land use for the 10 years before the start of the bare fallow experiment. The field was tilled twice a year, in autumn (October) and spring (April), to 0.25 m depth. Weeds were removed by hand and through the use of herbicides. SOC measurements were obtained from six replicates over 79 years of observations. |
| References | | | | | |
| Christensen (1990); Christensen and Johnston (1997); Bruun et al. (2003); Barré et al. (2010) | Morel et al. (1984); Houot et al. (1989); Barré et al. (2010) | Lazarev (2007); Barré et al. (2010); Guenet et al. (2013) | Christensen et al. (2006, 2019); Barré et al. (2010) | Kirchmann et al. (1994); Andrén and Katterer (1997); Gerzabek et al. (1997); Kirchmann and Gerzabek (1999); Barré et al. (2010); Kätterer et al. (2011) | Burgevin and Hénin (1939); Pernes-Debuyser and Tessier (2002); Barré et al. (2010); Paradelo et al. (2013, 2015); van Oort et al. (2018) |

Fig. A. Classification of sites with respect to De Martonne-Gottmann aridity index (De Martonne, 1942) and heat wave days' frequency. For the aridity index ($b$), the following range limits discriminate among thermo-pluviometric conditions associated with aridity gradients: $b<5$: extreme aridity; $5 \leq b \leq 14$: aridity; $15 \leq b \leq 19$: semi-aridity; $20 \leq b \leq 29$: sub-humidity; $30 \leq b \leq 59$: humidity; $b>59$: strong humidity. For identifying the frequency of heat wave ($hw$) days within a year in each site, we summarized the number of consecutive days (at least seven) when the maximum temperature was higher than the average summer (June, July and August) maximum temperature of all the available years (baseline) +3 °C (Confalonieri et al., 2010 after Barnett et al., 2006). The range limits in this study were given after Sándor et al. (2017): $hw \leq 14$: extremely moderate frequency; $14 < hw \leq 28$: very moderate frequency; $28 < hw \leq 42$: moderate frequency; $42 < hw \leq 56$: high frequency; $56 < hw \leq 70$: very high frequency; $hw>70$: extremely high frequency.

Fig. B1. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M01 for blind simulations (Blind) and Specific calibration scenario (as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
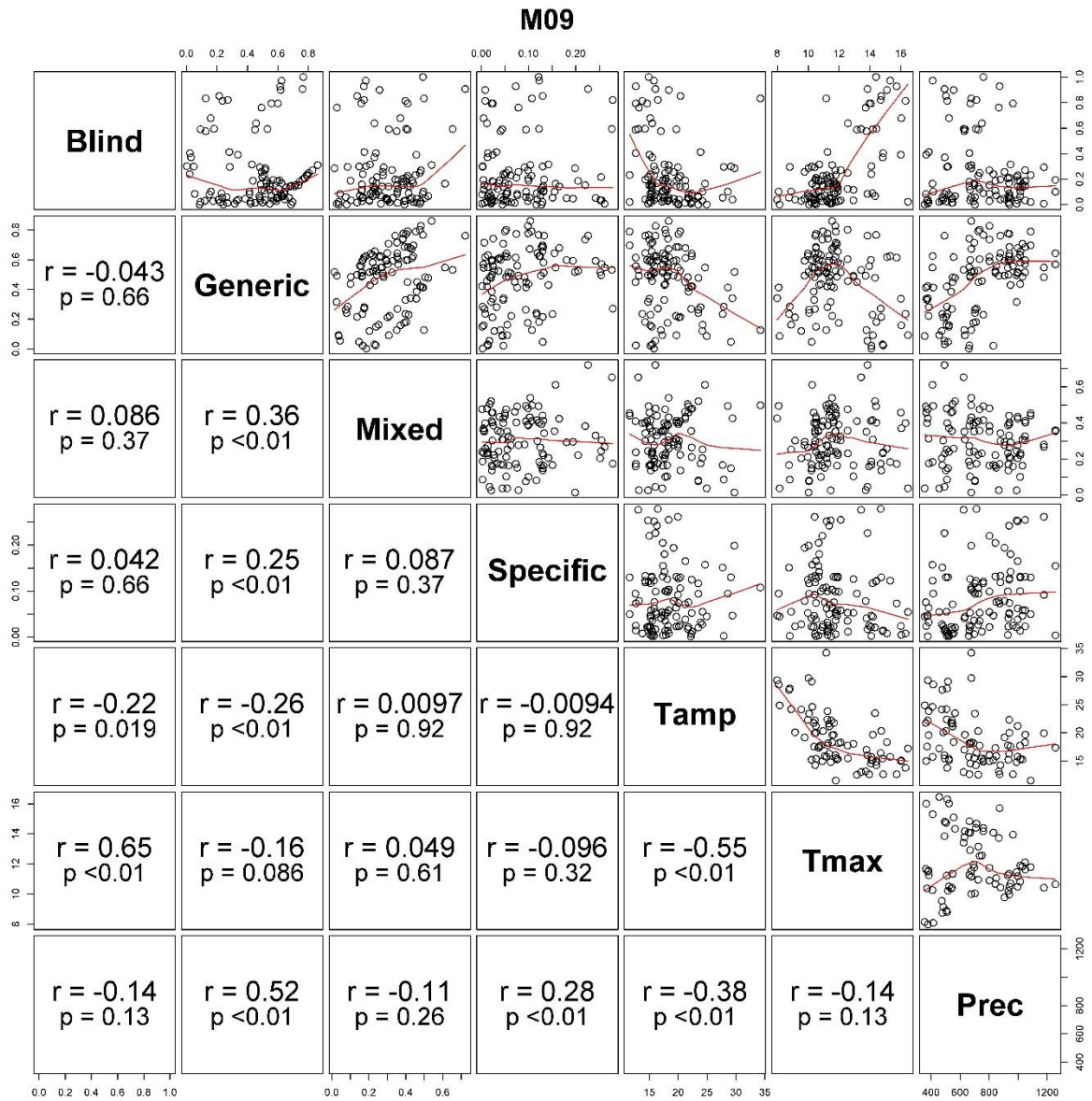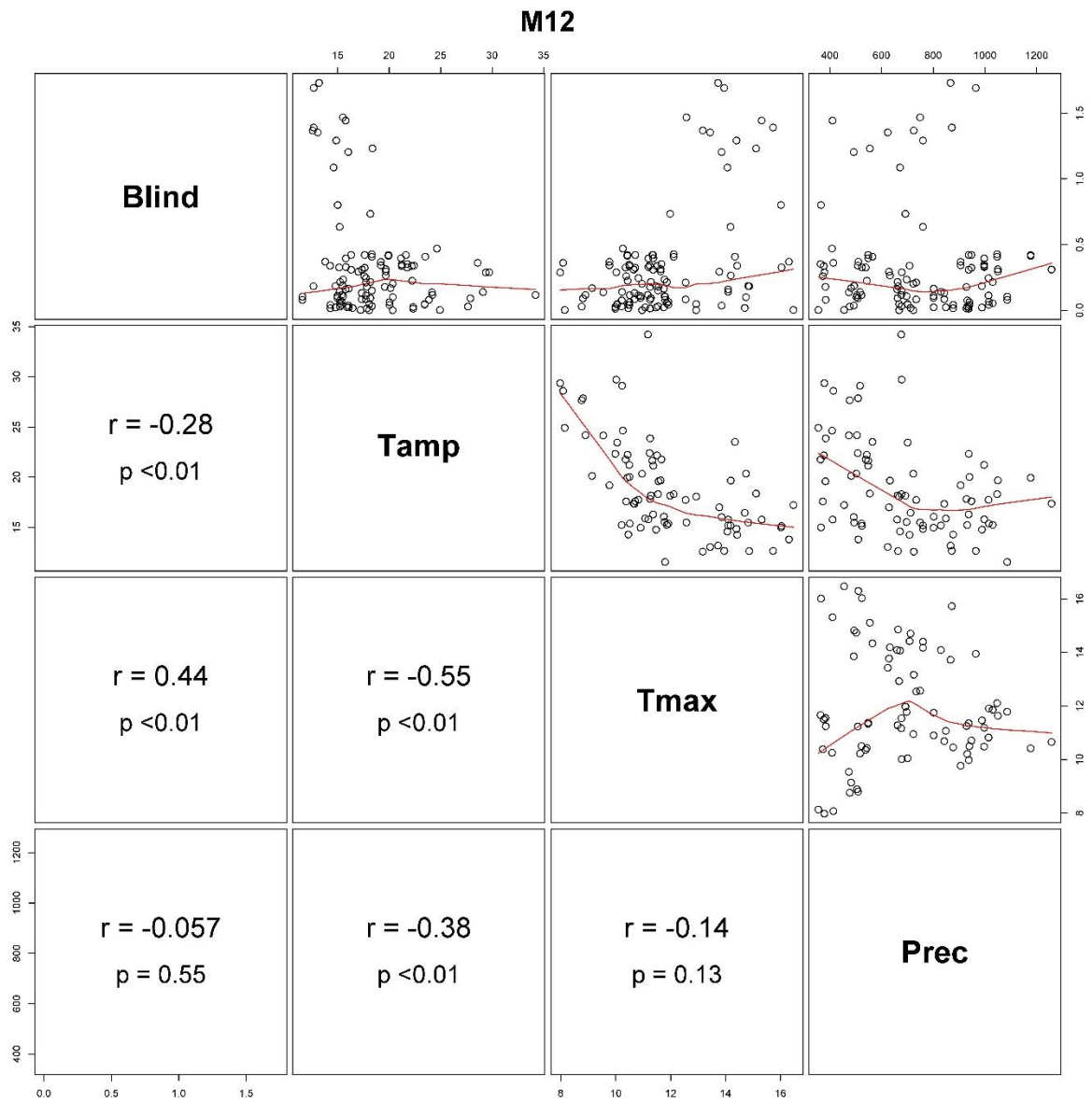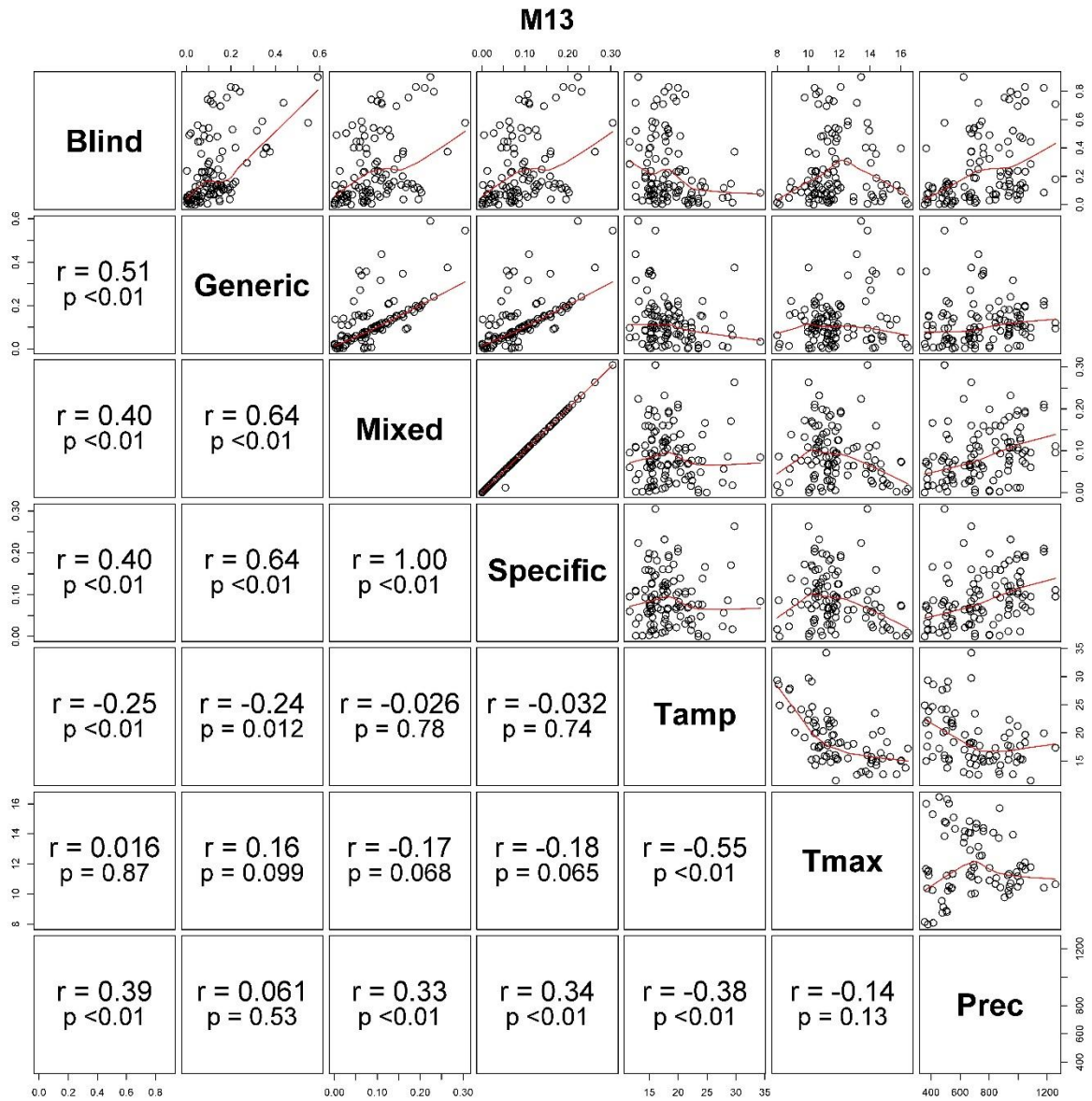
Fig. B2. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M02 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B3. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M03 for blind simulations (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
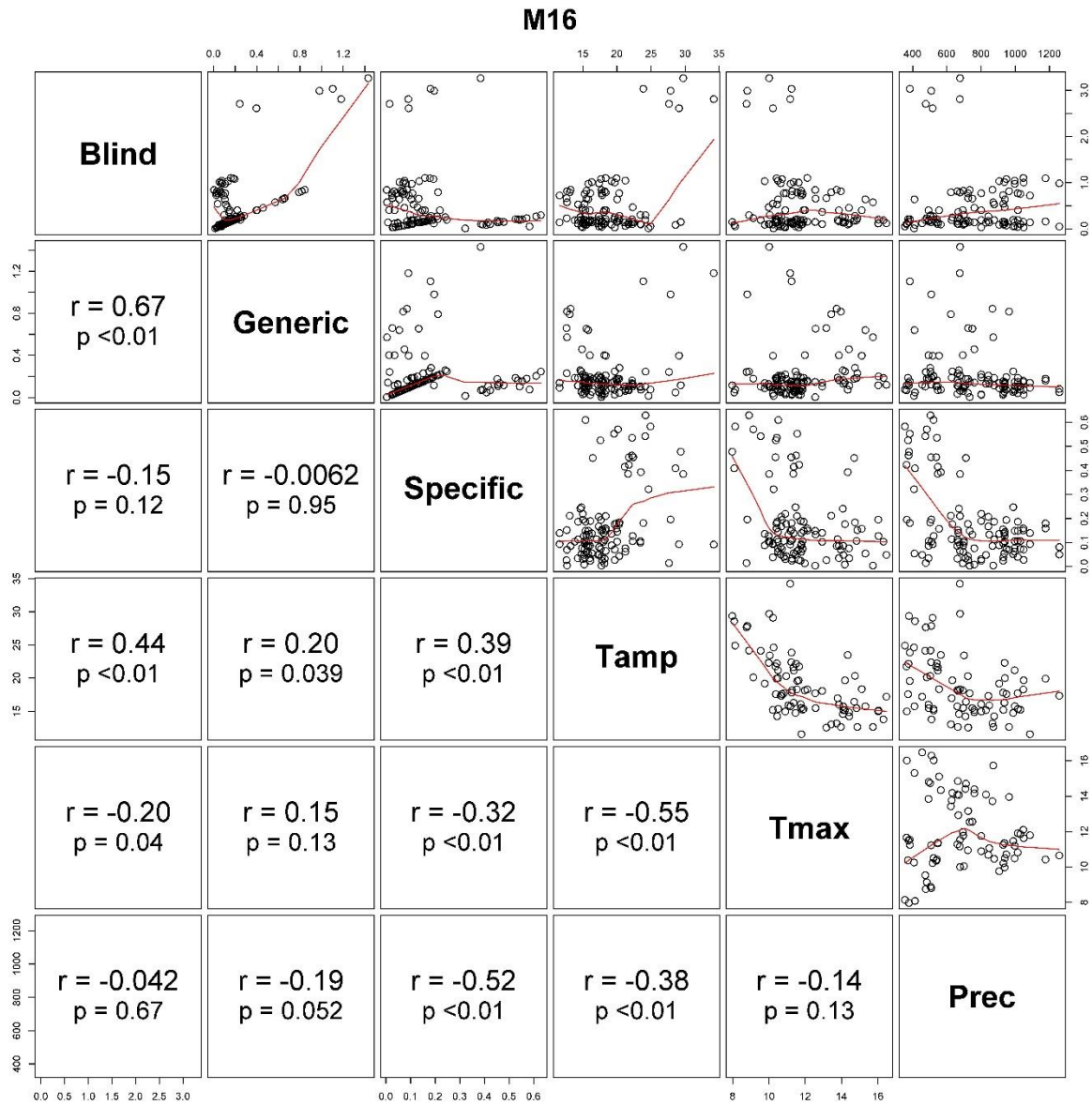
Fig. B4. Scatterplot correlation matrix of SOC (Mg C ha⁻¹) model residuals of M04 for blind simulations (Blind) and calibration scenarios (Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve. (Prec). Overlaid (red line) is a local non-parametric smoother curve.
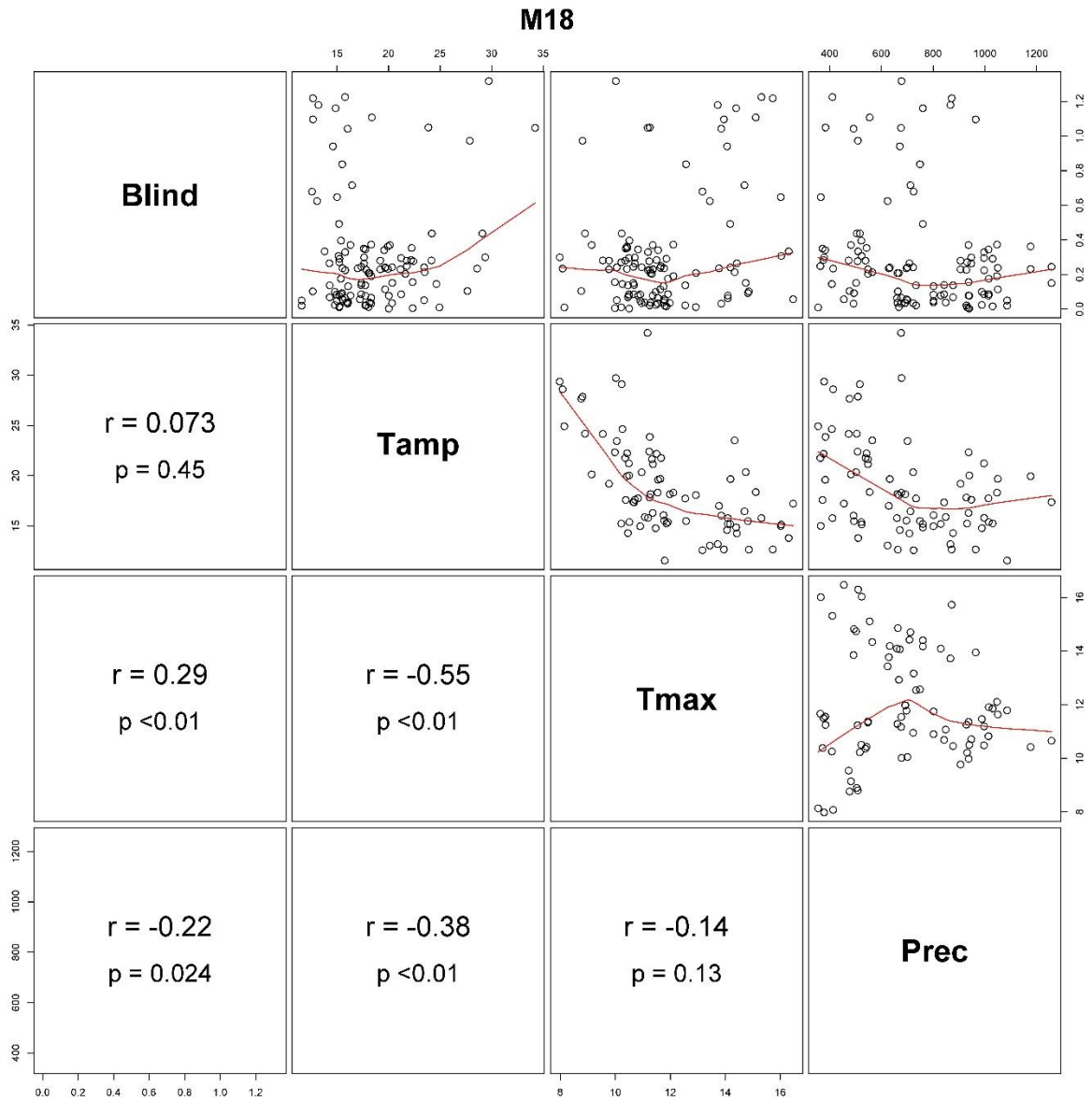
# M05



Fig. B5. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M05for blind simulations (Blind) and calibration scenarios (Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B6. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M06 for blind simulations (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
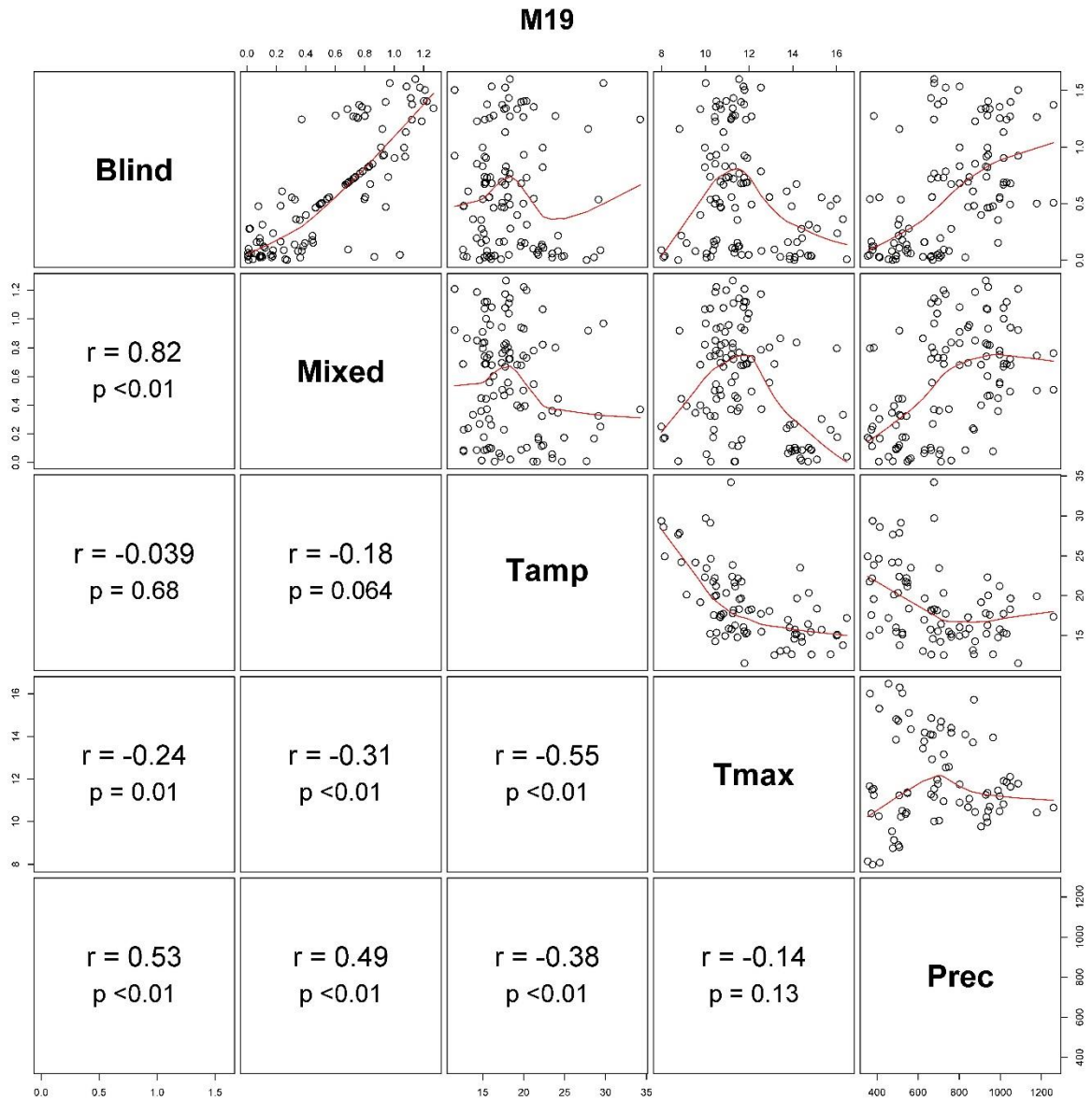
Fig. B7. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M07 for blind simulations (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
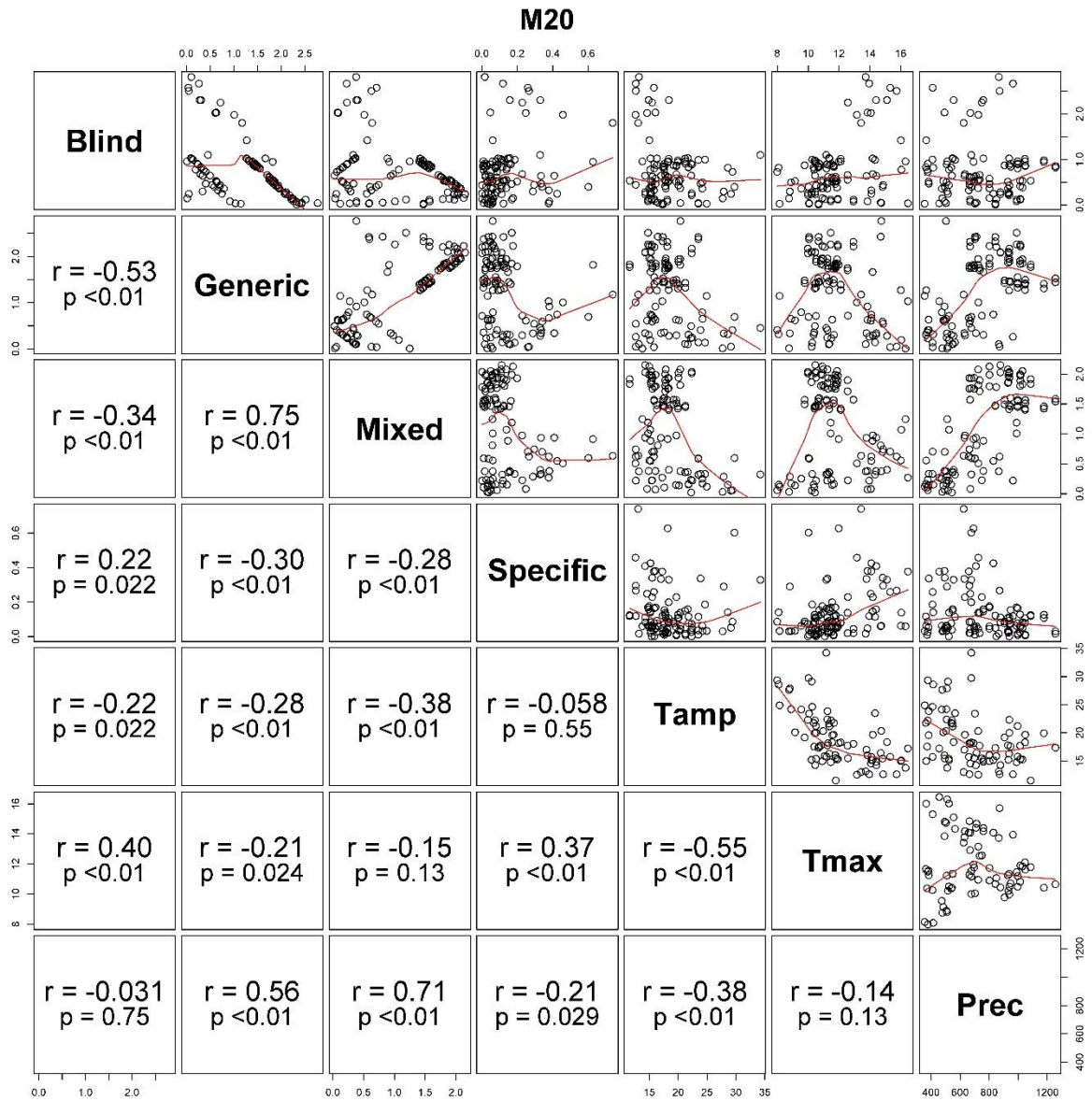
Fig. B8. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M09 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specifics in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B9. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M12 for blind simulations (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
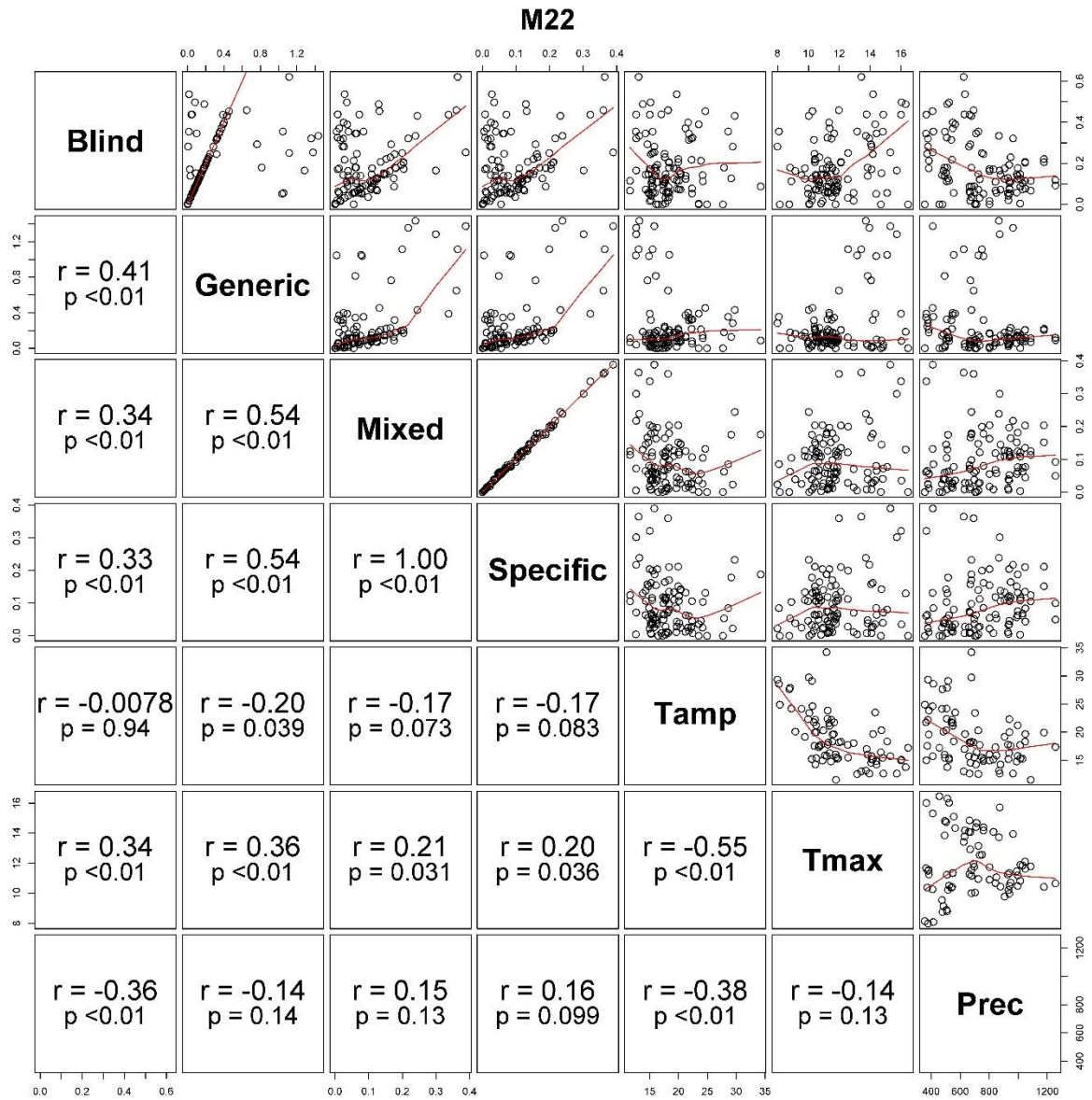
Fig. B10. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M13 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
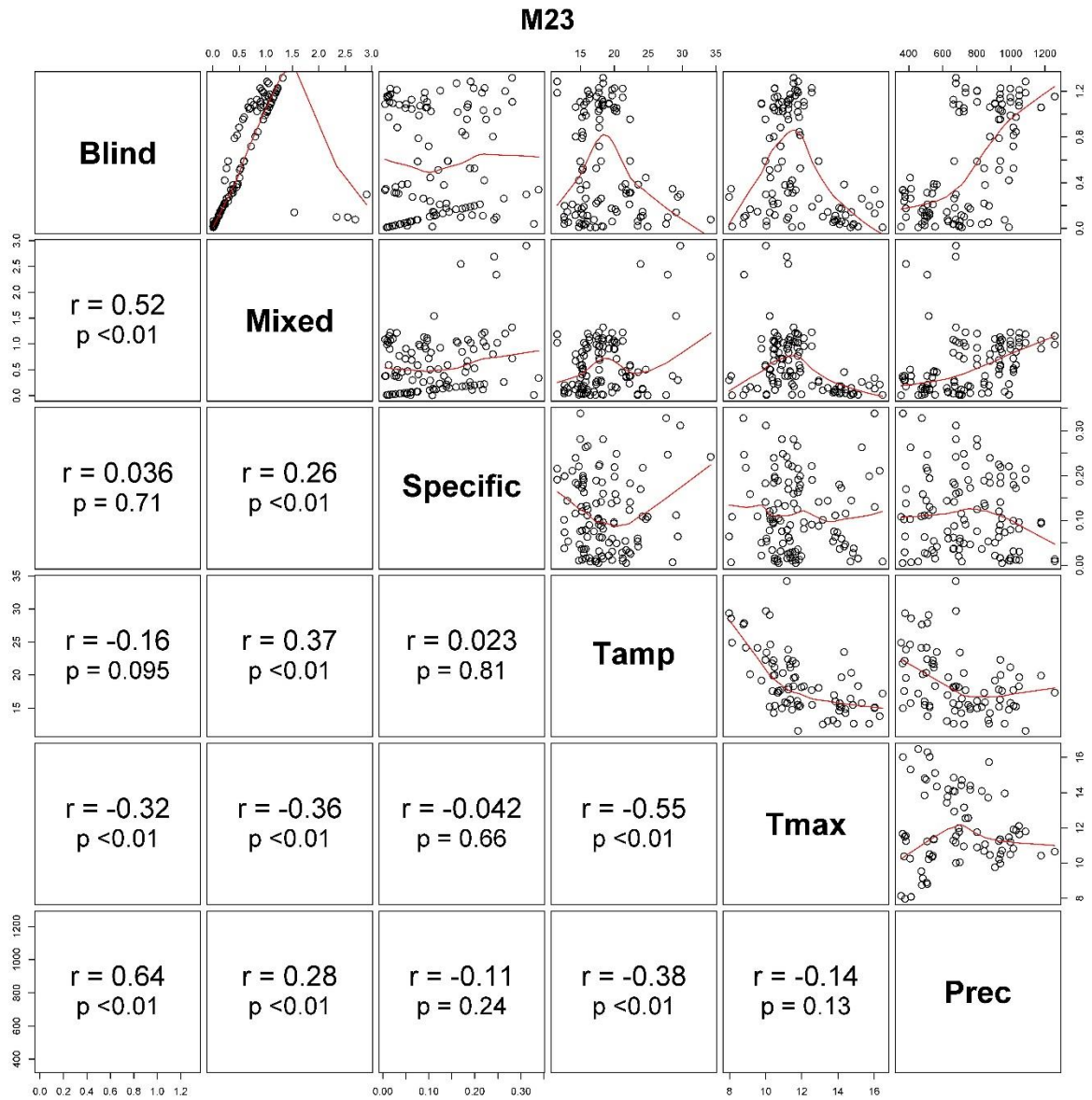
Fig. B11. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M16 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B12. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M18 for blind simulations (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
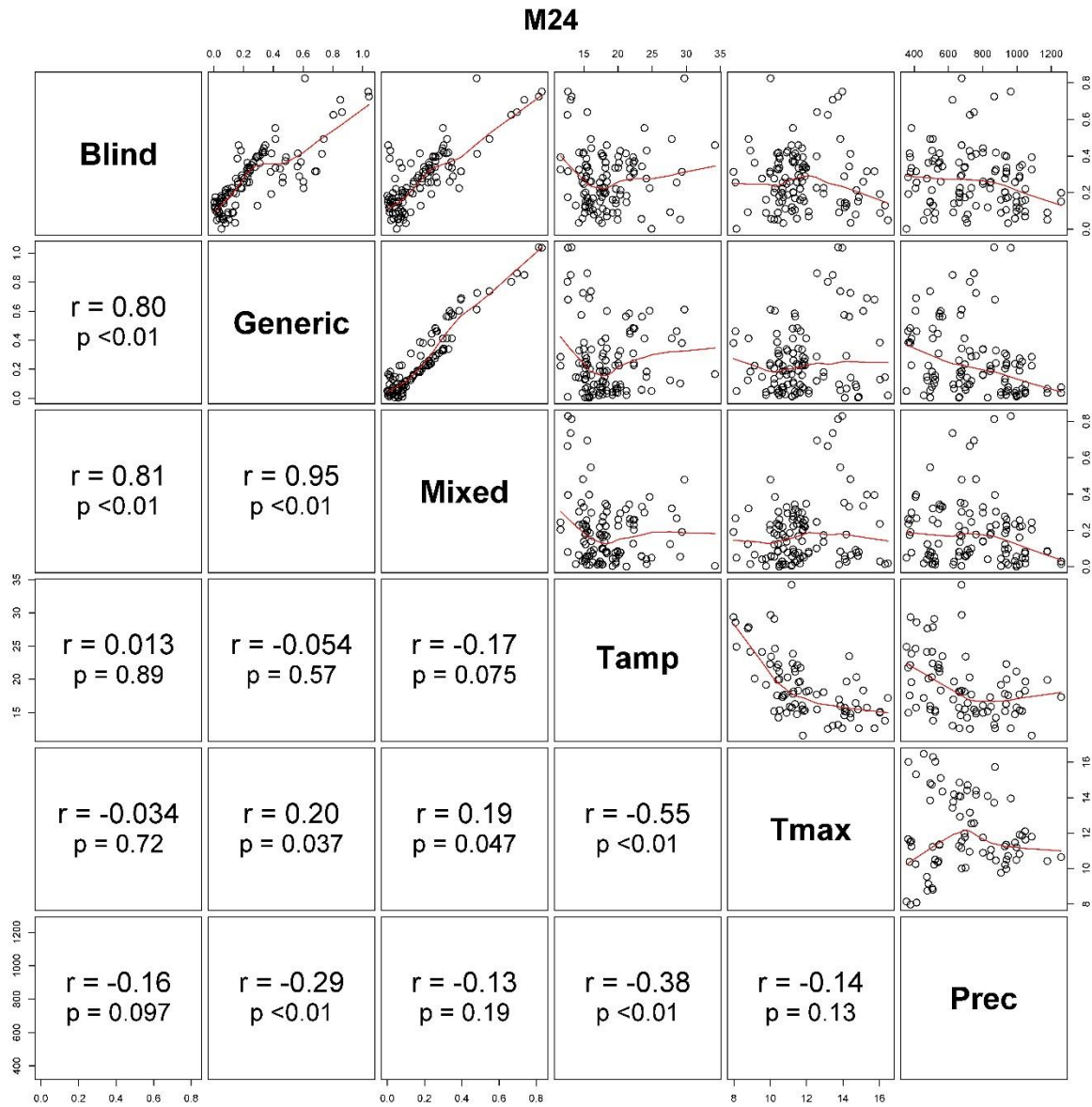
Fig. B13. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M19 for blind simulations (Blind) and Mixed calibration scenario (as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
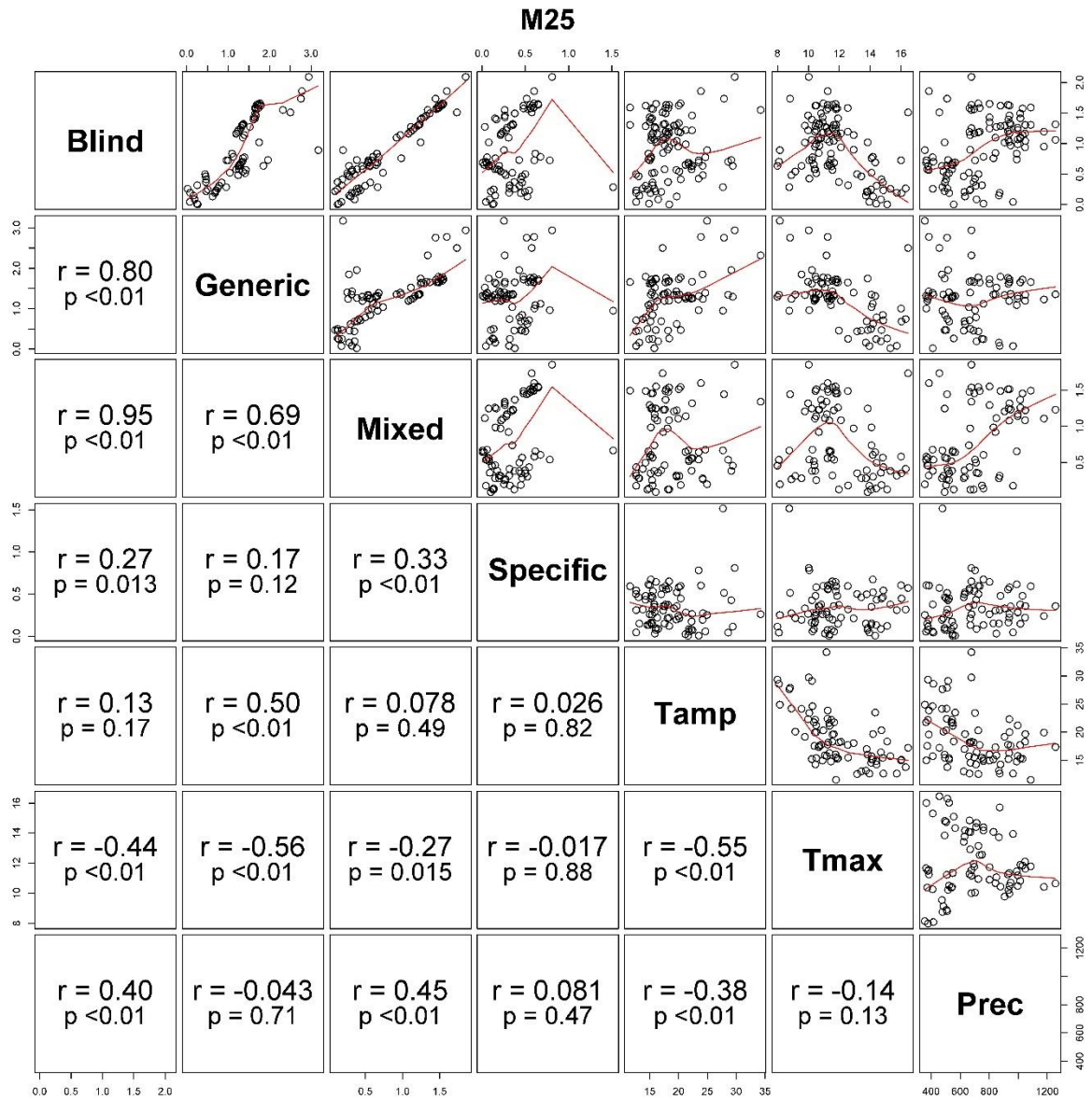
Fig. B14. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M20 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B15. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M22 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
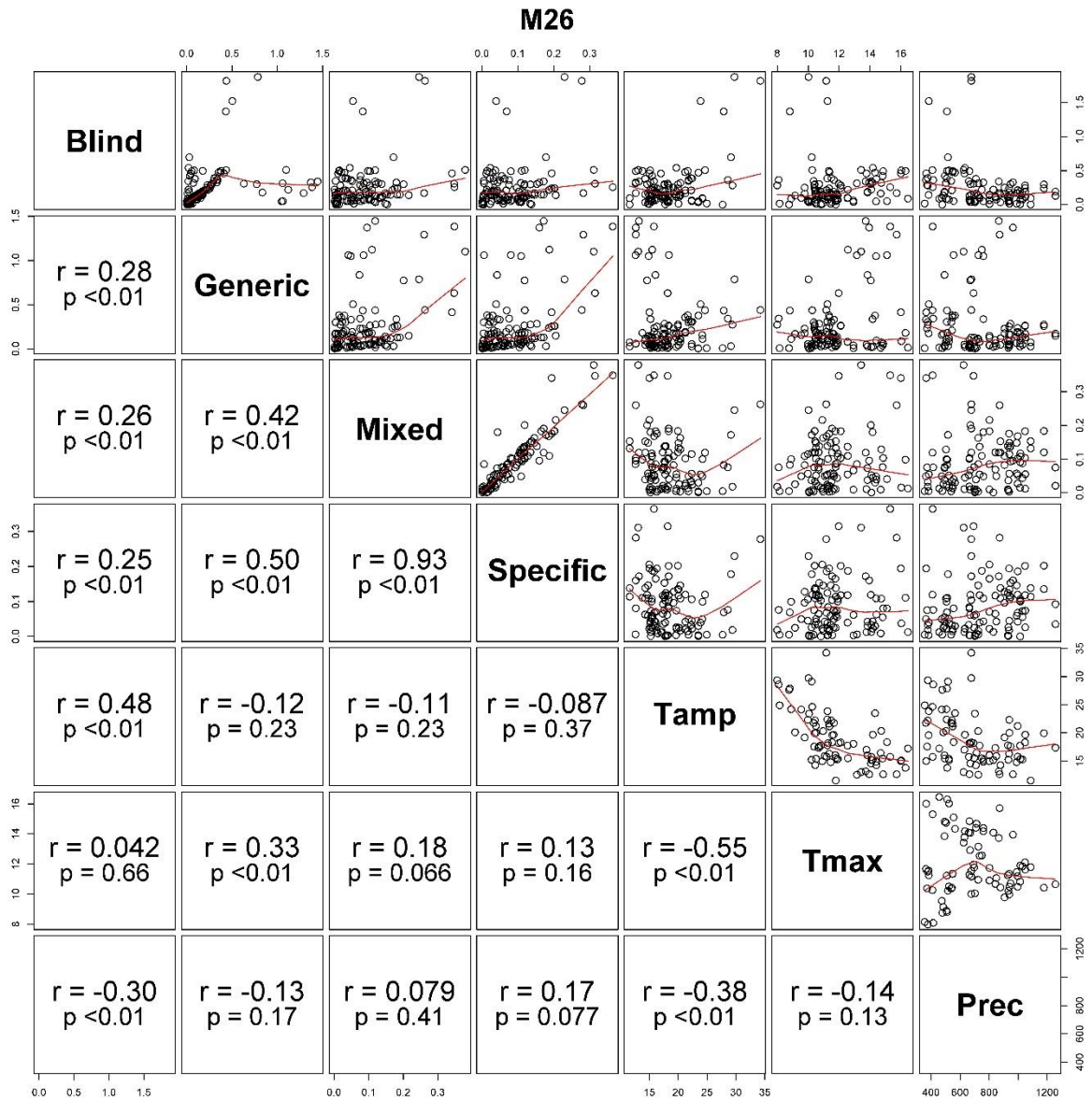
Fig. B16. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M23 for blind simulations (Blind) and calibration scenarios (Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
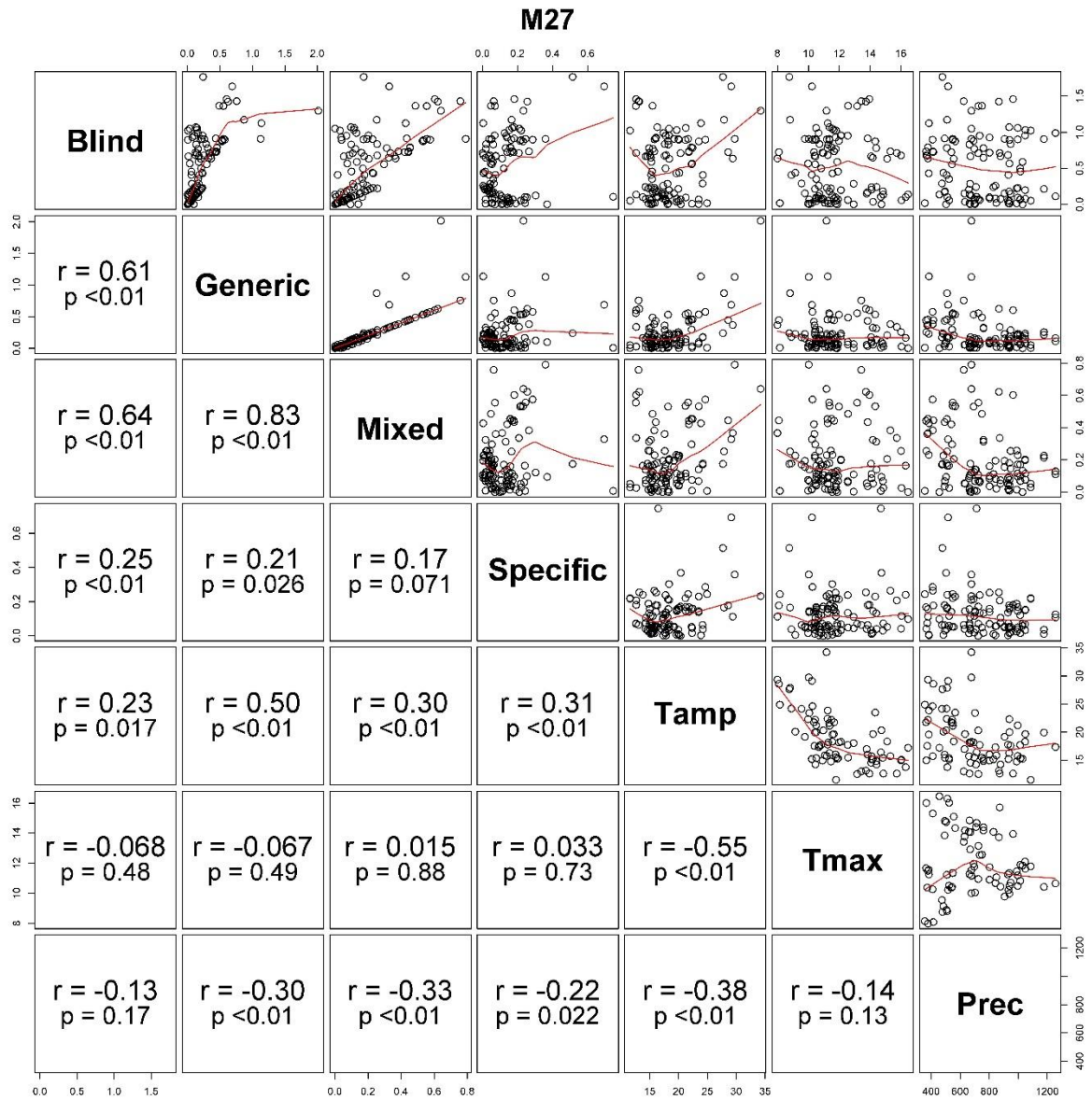
Fig. B17. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M24 for blind simulations (Blind) and calibration scenarios (Generic and Mixed as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B18. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M25 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
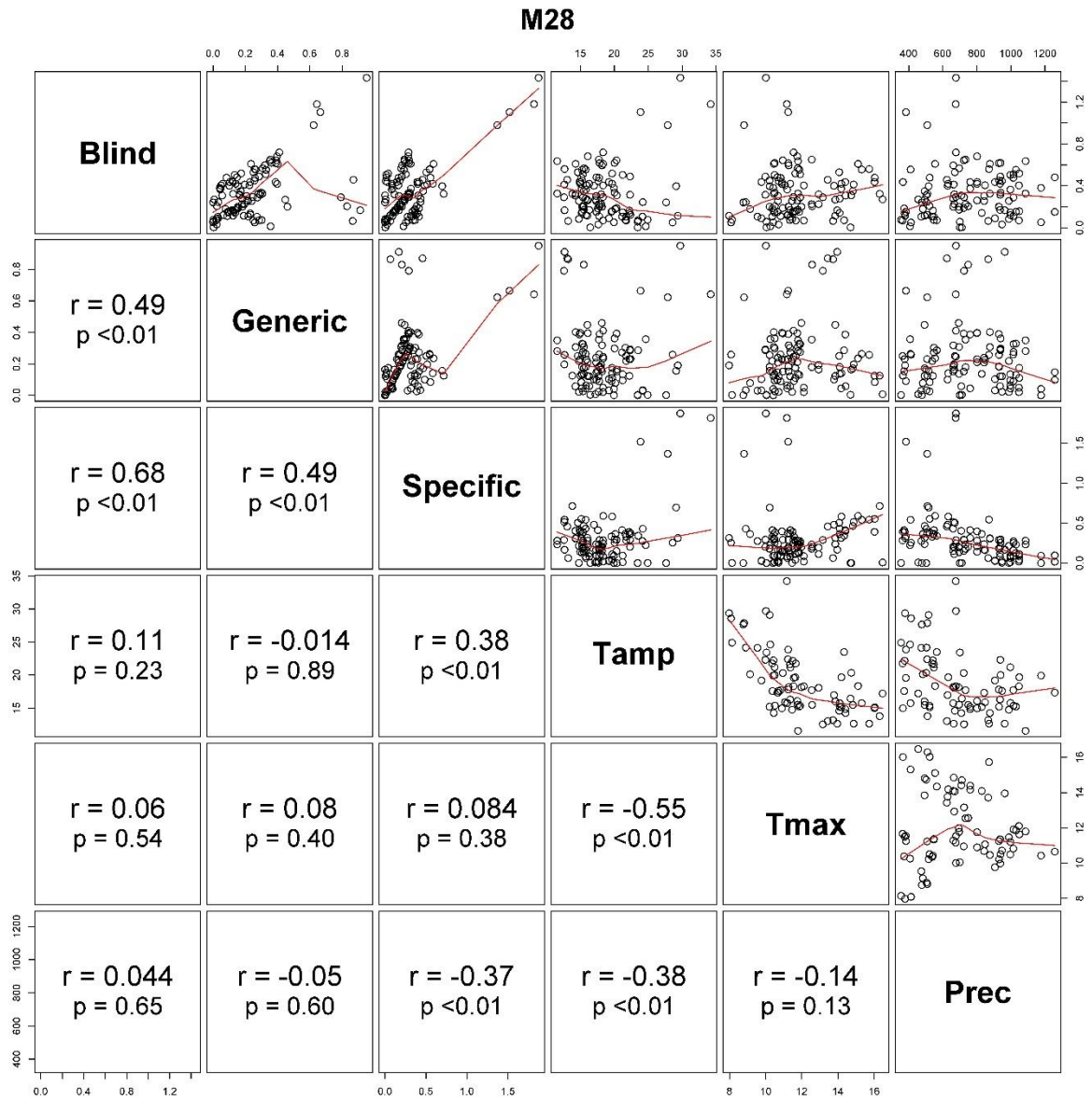
**M26**



Fig. B19. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M26 for blind simulations (Blind) and calibration scenarios (Generic, Mixed, Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
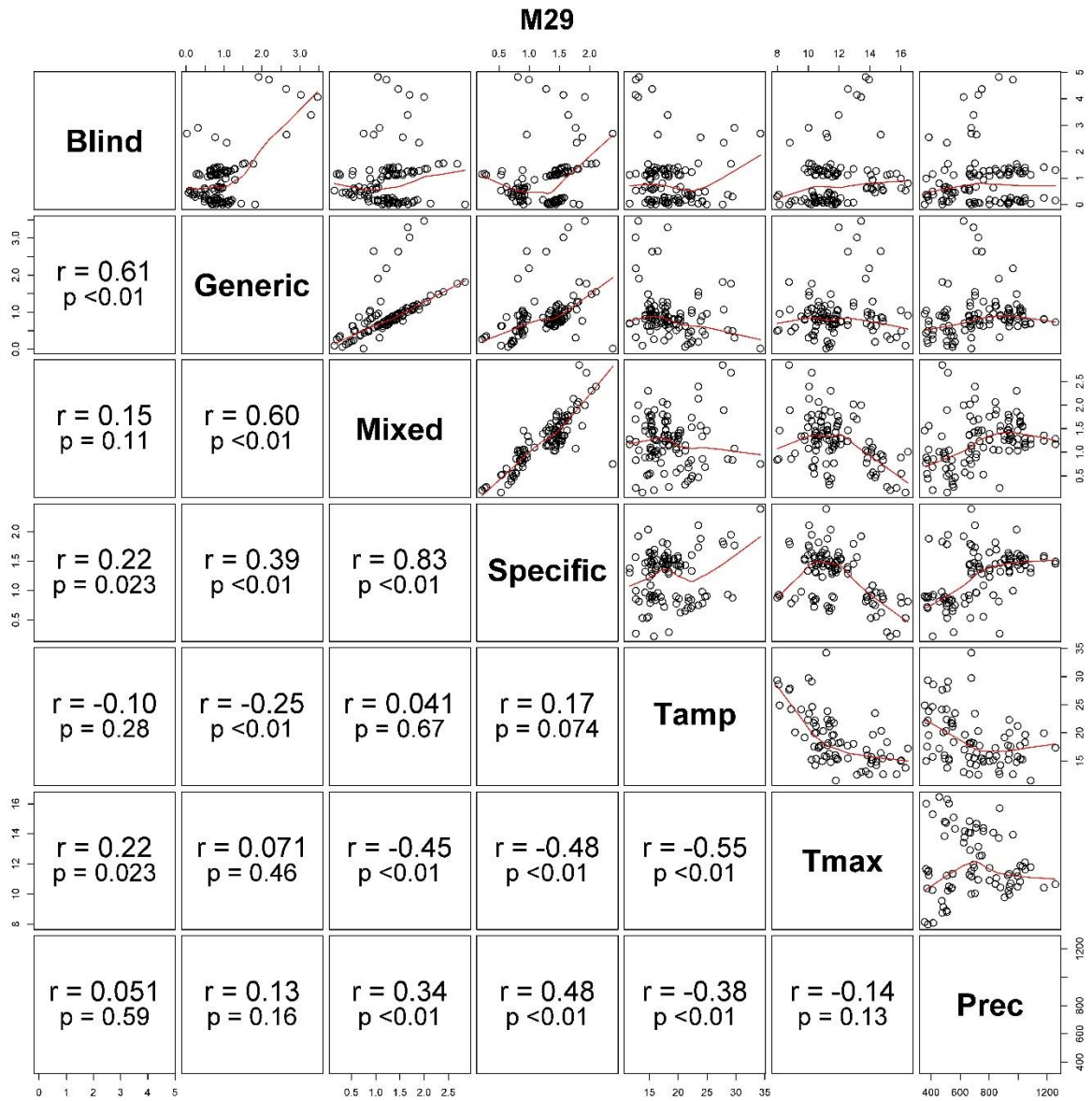
Fig. B20. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M27 for blind simulations (Blind) and calibration scenarios (Generic, Mixed, and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B21. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M28 for blind simulations (Blind) and calibration scenarios (Generic and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
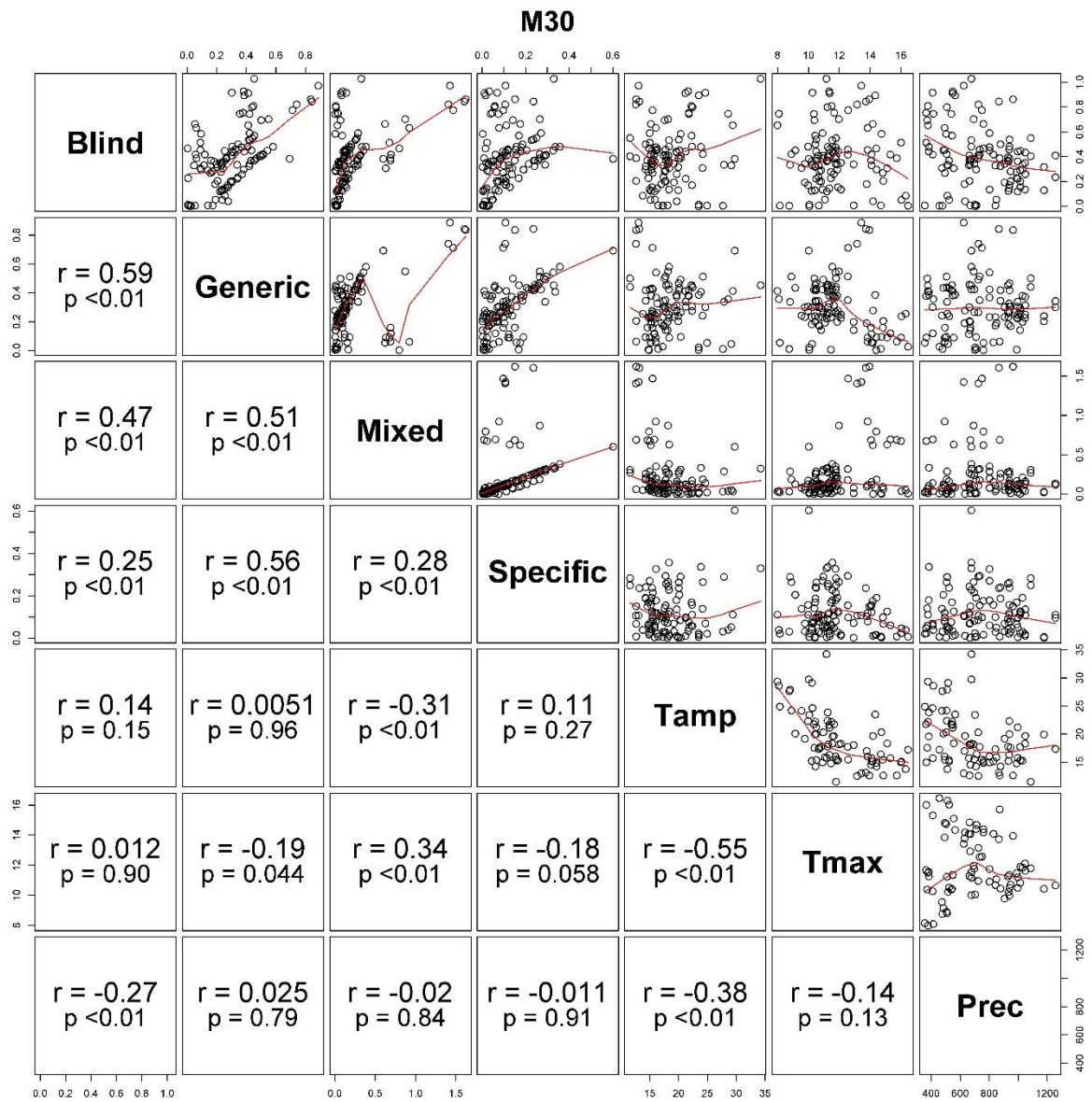
Fig. B22. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M29 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
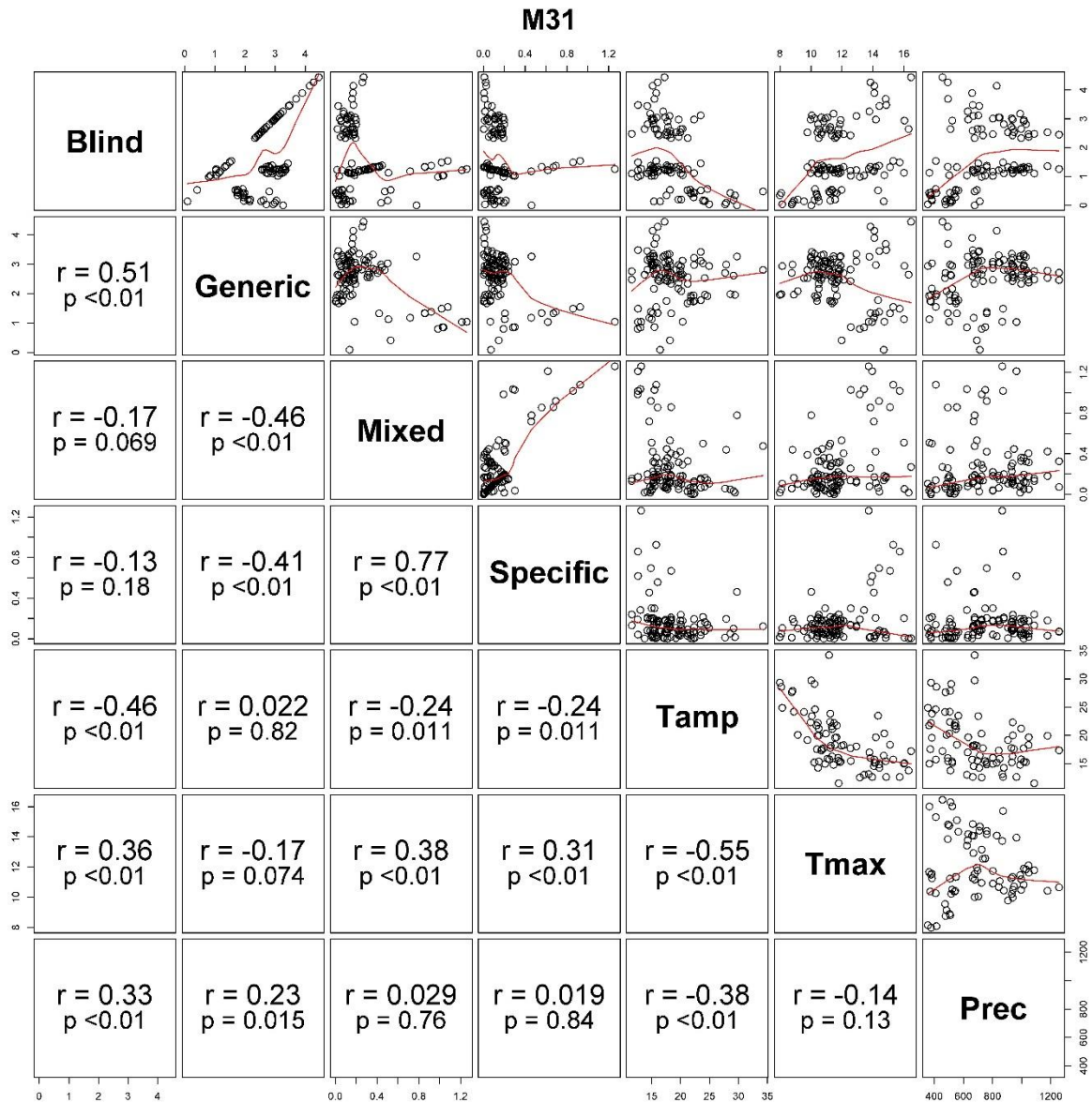
Fig. B23. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M30 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
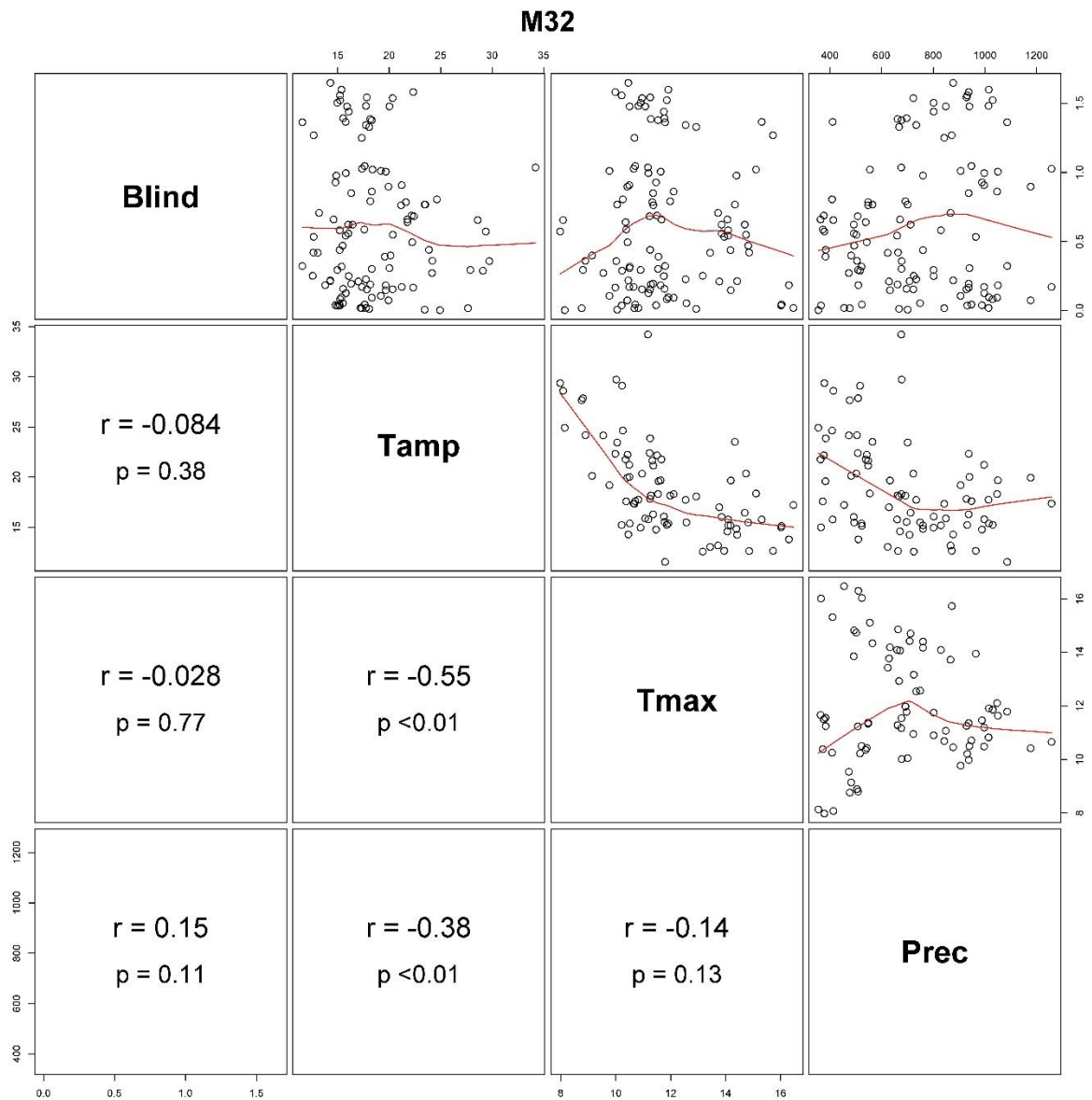
Fig. B24. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M31 for blind simulations (Blind) and calibration scenarios (Generic, Mixed and Specific as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

Fig. B25. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M32 for blind (Blind) and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.
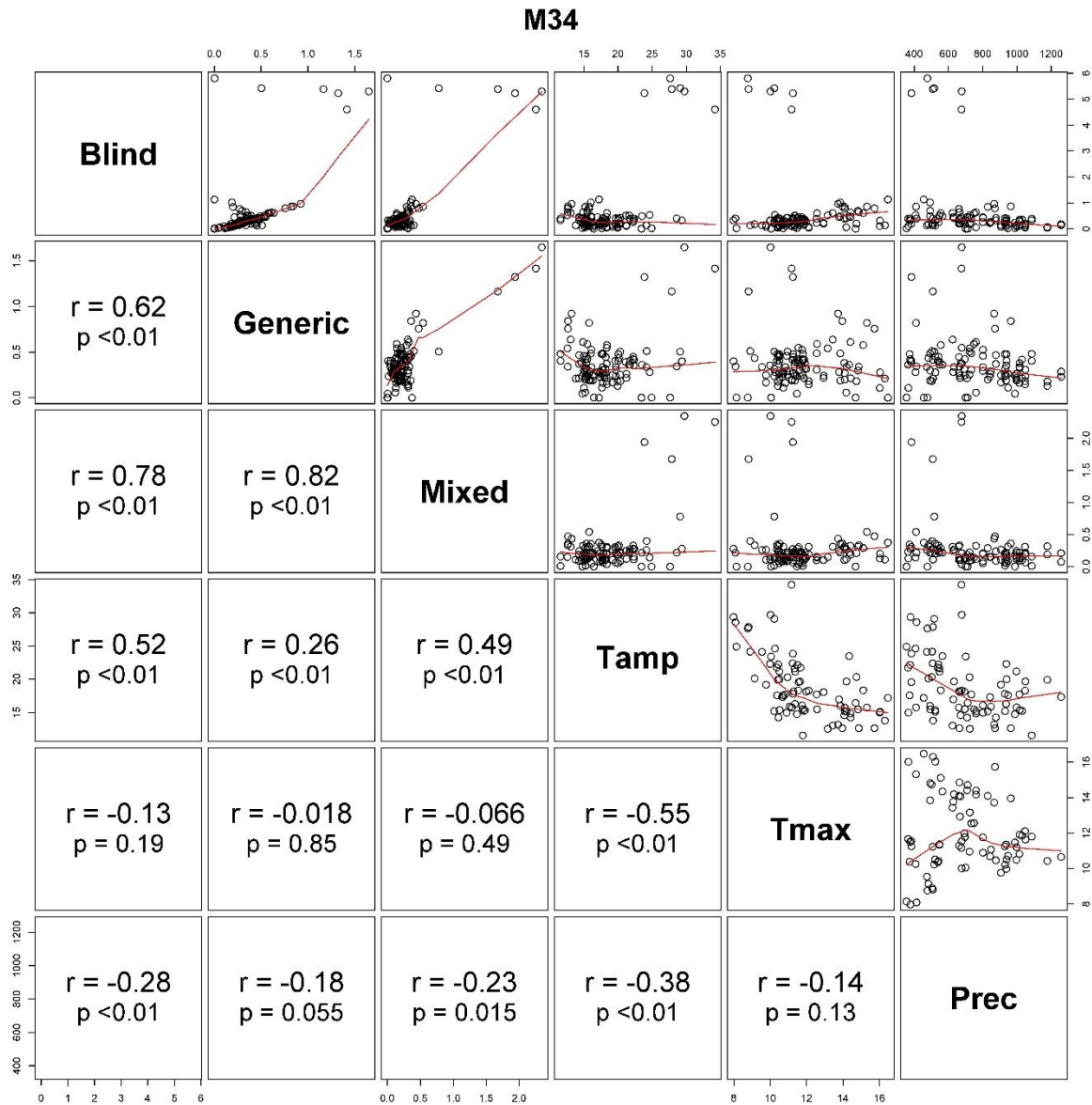
Fig. B26. Scatterplot correlation matrix of SOC (Mg C ha$^{-1}$) model residuals of M34 for blind simulations (Blind) and calibration scenarios (Generic and Mixed as in Table 3), and the annual climate metrics maximum temperature (Tmax), mean temperature amplitude (Tamp) and precipitation (Prec). Overlaid (red line) is a local non-parametric smoother curve.

**References**

Burgevin, H., & Hénin, S. (1939). Dix années d'expériences sur l'action des engrais sur la composition et les propriétés d'un sol de limon. *Annales Agronomiques*, **9**, 771-799. (in French)

Andrén, O., & Kätterer, T. (1997). ICBM: The introductory carbon balance model for exploration of soil carbon balances. *Ecological Applications*, **7**, 1226-1236. https://doi.org/10.1890/1051-0761(1997)007[1226:ITICBM]2.0.CO;2

Barnett, C., Hossel, J., Perry, M., Procter, C., & Hughes, G. (2006). A handbook of climate trends across Scotland. Edinburgh (UK): Scotland and Northern Ireland Forum for Environmental Research, SNIFFER Project CC03.

Bruun, S., Christensen, B. T., Hansen, E. M., Magid, J., & Jensen, L. S. (2003). Calibration and validation of the soil organic matter dynamics of the Daisy model with data from the Askov long-term experiments. *Soil Biology and Biochemistry*, **35**, 67-76. https://doi.org/10.1016/S0038-0717(02)00237-7

Christensen, B. T. (1990). Effect of cropping system on the soil organic matter content II, Field experiments on a sandy loam 1956-1985. *Tidsskrift for Planteavl*, **94**, 161-169. (in Danish with English summary)

Christensen, B. T., & Johnston, A. E. (1997). Soil organic matter and soil quality - lessons learned from long-term experiments at Askov and Rothamsted. *Developments in Soil Science*, **25**, 399-430. https://doi.org/10.1016/S0166-2481(97)80045-1

Christensen, B. T., Petersen, J., & Trentemoller, U. M. (2006) The Askov long-term experiments on animal manure and mineral fertilizers: The Lermarken site 1894-2004, DIAS Report Plant Production no. 121. Tjele (Denmark): Danish Institute of Agricultural Sciences.

Christensen, B. T., Thomsen, I. K., & Eriksen, J. (2019). The Askov long-term experiments: 1894-2019. A unique research platform turns 125 years. DCA report, No. 151. Aarhus (Denmark): Aarhus University, Danish Centre for Food and Agriculture.

Confalonieri, R., Bellocchi, G., & Donatelli, M. (2010). A software component to compute agro-meteorological indicators. *Environmental Modelling & Software*, **25**, 1485-1486. https://doi.org/10.1016/j.envsoft.2008.11.007

De Martonne, E. (1942). Nouvelle carte mondiale de l'indice d'aridité. *Annales de Géographie*, **51**, 242-250. (in French)

Gerzabek, M. H., Pichlmayer, F., Kirchmann, H., & Haberhauer, G. (1997). The response of soil organic matter to manure amendments in a long-term experiment at Ultuna, Sweden.

*European Journal of Soil Science*, **48**, 273-282. https://doi.org/10.1111/j.1365-2389.1997.tb00547.x

Guenet, B., Eglin, T., Vasilyeva, N., Peylin, P., Ciais, P., & Chenu, C. (2013). The relative importance of decomposition and transport mechanisms for soil organic carbon profiles. *Biogeosciences*, **10**, 2379-2392. https://doi.org/10.5194/bg-10-2379-2013

Houot, S., Molina, J. A. E., Chaussod, R., & Clapp, C. E. (1989). Simulation by NCSOIL of net mineralization in soils from the Deherain and 36 parcelles fields at Grignon. *Soil Science Society of America Journal*, **53**, 451-455. Retrieved from https://acsess.onlinelibrary.wiley.com/doi/pdf/10.2136/sssaj1989.03615995005300020023x

Kätterer, T., Bolinder, M. A., Andrén, O., Kirchmann, H., & Menichetti, L. (2011). Roots contribute more to refractory soil organic matter than aboveground crop residues, as revealed by a long-term field experiment. *Agriculture, Ecosystems & Environment*, **141**, 184-192. https://doi.org/10.1016/j.agee.2011.02.029

Kirchmann, H., & Gerzabek, M. H. (1999). Relationship between soil organic matter and micropores in a long-term experiment at Ultuna, Sweden. *Journal of Plant Nutrition and Soil Science*, **162**, 493-498. https://doi.org/10.1002/(SICI)1522-2624(199910)162:5<493::AID-JPLN493>3.0.CO;2-S

Kirchmann, H., Persson, J., & Carlgren, K. (1994). The Ultuna long-term soil organic matter experiment, 1956-1991. Department of Soil Sciences, Reports and Dissertations, 17. Uppsala (Sweden): Swedish University of Agricultural Sciences.

Lazarev, V. I. (2007). Dynamics of agro-physical soil properties. In *Dynamics of effective fertility of chernozem under long-term agricultural use* (pp. 89-94). Kursk (Russia): Kursk State Agricultural Academy. (in Russian)

Morel, R., Lasnier, T., & Bourgeois, P. (1984). Les essais de fertilisation de longue durée de la station agronomique de Grignon ; Dispositif Dehérain et des 36 parcelles : Résultats expérimentaux (période 1938-1982). Paris (France): Institut National de la Recherche Agronomique. (in French)

Paradelo, R., van Oort, F., & Chenu, C. (2013). Water-dispersible clay in bare fallow soils after 80 years of continuous fertilizer addition. *Geoderma*, **200-201**, 40-44.

Paradelo, R., Virto, I., & Chenu, C. (2015). Net effect of liming on soil organic carbon stocks: A review. *Agriculture, Ecosystems & Environment*, **202**, 98-107. https://doi.org/10.1016/j.agee.2015.01.005

Pernes-Debuyser, A., & Tessier, D. (2002). Influence du pH sur les propriétés des sols : l'essai de longue durée des 42 parcelles à Versailles. *Revue de Sciences de l'Eau*, **15**, 27-39. (in French with English summary)

Sándor, R., Barcza, Z., Acutis, M., Doro, L., Hidy, D., Köchy, M., … Bellocchi, G. (2017). Multi-model simulation of soil temperature, soil water content and biomass in Euro-Mediterranean grasslands: Uncertainties and ensemble performance. *European Journal of Agronomy*, **88**, 22-40. https://doi.org/10.1016/j.eja.2016.06.006

van Oort, F., Paradelo, R., Proix, N., Delarue, G., Baize, D., & Monna, F. (2018). Centennial fertilization-induced soil processes control trace metal dynamics. Lessons from a long-term bare fallow experiment. *Soil Systems*, 2, 23. https://doi.org/10.3390/soilsystems2020023