1    **Molecular tradeoffs in soil organic carbon composition at continental scale**

2

3    Steven J. Hall[1]*†, Chenglong Ye[1,2]*, Samantha R. Weintraub[3], William C. Hockaday[4]†

4

5    [1]Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, IA

6    50011, USA

7    [2]Ecosystem Ecology Lab, College of Resources and Environmental Sciences, Nanjing

8    Agricultural University, Nanjing 210095, China

9    [3]National Ecological Observatory Network, Battelle, Boulder, Colorado, 80301

10    [4]Department of Geosciences, Baylor University, Waco, TX 76798, USA

11    *Contributed equally to this manuscript

12    †Co-corresponding authors: stevenjh@iastate.edu; Phone: 515-294-7650; Fax: 515-294-1337

13    william_hockaday@baylor.edu; Phone: 254-710-2639; Fax: 254-710-2673

14

15

16

17

18

19

20

21

22

23

24

**The molecular composition of soil organic carbon remains contentious. Microbial-, plant-, and fire-derived compounds may each contribute, but do they vary predictably among ecosystems? Here we present carbon functional groups and molecules from a diverse spectrum of North American surface mineral soils, primarily collected from the National Ecological Observatory Network, quantified by nuclear magnetic resonance spectroscopy and a molecular mixing model. Soils varied widely in relative contributions of carbohydrate, lipid, protein, lignin, and char-like carbon, but each compound class had similar overall abundance. Three principal component axes explained 90% of the variance in carbon composition: the first showed a tradeoff between lignin and protein, the second showed a tradeoff between carbohydrate and char, and the third was explained by lipids. Reactive aluminum, crystalline iron oxides, and pH plus overlying organic horizon thickness best explained variation along each respective axis; these predictors were ultimately related to climate. Together, our data point to continental-scale tradeoffs in soil carbon molecular composition which are linked to environmental and geochemical variables known to predict carbon mass concentrations. Controversies regarding the genesis of soil carbon and its potential responses to global change can be partially reconciled by considering diverse ecosystem properties that drive complementary persistence mechanisms.**

Soil organic carbon (SOC) is generally understood to comprise a diverse suite of biomolecules representing the decomposition products of plant and microbial biomass and the imprint of abiotic processes such as fire[1,2]. However, the fundamental mechanisms controlling the molecular composition of SOC within and among mineral soils remain contentious[3]. Do disparate soils converge along a predictable molecular continuum of SOC composition driven by

48 the inexorable transformation of plant detritus to a consistent suite of low-molecular-weight

49 decomposition products[1,4–6]? Or conversely, do diverse biogeochemical factors such as climate,

50 vegetation, or mineralogy lead to distinct molecular differences in SOC among ecosystems[7,8]?

51 We can increasingly predict the spatial distribution of SOC as a function of biogeochemical

52 properties[9], as well as the partitioning of SOC between particulate and mineral-associated

53 pools[10]. An equivalent framework for predicting SOC molecular composition among ecosystems

54 remains elusive but could inform our understanding of the functional properties of SOC and its

55 dynamics under global change[3,11].

56     Debates on the importance of different mechanisms of SOC persistence rest in part on our

57 contested understanding of its molecular composition. Plant-derived aromatic compounds like

58 lignin were historically thought to dominate SOC due to their macromolecular structure[12], which

59 requires strong oxidants for depolymerization[13]. Subsequent work challenged this view by

60 demonstrating that aromatic and lignin-like moieties may be minor constituents[5] that decompose

61 faster than bulk SOC[14,15]. Microbial necromass and low-molecular-weight decomposition

62 products (carbohydrates, proteins, and lipids) have assumed key roles in current SOC paradigms

63 given the potential for efficient microbial metabolism and recycling of these molecules. In this

64 view, SOC persistence does not derive from chemical complexity[16] or stability but rather from

65 protective physico-chemical interactions with minerals and aggregates with microbial detritus

66 playing a dominant role[6,17–21].

67     However, the importance of microbial vs. direct plant contributions to SOC could vary

68 among ecosystems[22]. Microbial growth and necromass production may be decoupled from SOC

69 accumulation in stressful environments where decomposition is inefficient[23]. Significant

70 contributions of lignin and other plant-derived compounds to mineral-associated SOC were

71    recently observed[24–26]. In fact, lignin-derived C may have been systematically underestimated

72    due to methodological biases[25,27,28]. Finally, char-like molecules presumably derived from

73    pyrolytic decomposition are prevalent in many ecosystems[29,30] despite evidence that the

74    increased stability of these molecules does not guarantee long-term persistence[31,32]. In spite of

75    significant theoretical and empirical progress[3,10,33], we still lack a consistent framework for

76    reconciling potential controls on SOC molecular composition across diverse ecosystem types.

77    Such a framework would enhance our conceptual understanding of the origins and persistence of

78    soil carbon to inform modeling and management of this critical resource.

79          Here, we leveraged a unique sample archive and datasets provided by the National

80    Ecological Observatory Network (NEON), along with additional samples, to characterize SOC

81    molecular composition and its relationships with biogeochemical factors across 42 North

82    American surface mineral soils (Extended Data Fig. 1). Samples spanned 11 of the 12 US

83    Department of Agriculture soil orders (all except Histosols, which were explicitly excluded) and

84    the major ecosystem gradients of North America (tropics to tundra; Supplementary Table 1,

85    Extended Data Fig. 2). We quantified the molecular functional groups of bulk SOC of

86    demineralized samples using solid-state $^{13}$C cross-polarization magic-angle-spinning (CPMAS)

87    nuclear magnetic resonance (NMR) spectroscopy. We confirmed the robustness of results by

88    assessing sample pretreatment and analytical biases in companion measurements employing

89    cross- and direct-polarization NMR (Supplementary Discussion; Supplementary Tables 1–6,

90    Supplementary Figures 1–5).

91

92    **Molecular variation in SOC at the continental scale**

93    The $^{13}$C CPMAS NMR spectra (Extended Data Fig. 3) illustrated substantial variability in SOC

94    molecular composition (Fig. 1a). Across all samples, O-alkyl, alkyl, and aromatic C were the

95    largest constituents, with mean values of 23%, 22%, and 21% of total SOC, respectively; each

96    varied as much as 2.5-, 3.4-, and 2.8-fold among samples (Fig. 1a). Amide/carboxyl C was less

97    abundant (13%; $P < 0.05$) but was greater than phenolic, N-alkyl/methoxyl, and di-O-alkyl C,

98    which each comprised 7% of SOC, on average (Fig. 1a). These results challenge a previous

99    literature synthesis of similar measurements, where disparate soils tended towards a consistent

100   ranking of C functional groups (O-alkyl > alkyl > aromatic > carbonyl)[5]. Across our diverse

101   soils, no constituent was dominant overall, and either O-alkyl, alkyl, or aromatic C could

102   predominate within an individual sample.

103           The molecular mixing model implied that five constituent molecules had similar relative

104   abundance across the dataset as a whole: carbohydrate, lignin, lipid, protein, and char comprised

105   mean values of 21%, 21%, 18%, 18%, and 17% of SOC, respectively (Fig. 1b). Despite their

106   similar means, these molecules varied greatly among soils, by as much as 4-, 33-, 46-, 10-, and

107   6-fold, respectively. Carbonyl C which represented the oxidized products of various molecules

108   was consistently less abundant (5% of SOC). To illustrate molecular changes during

109   decomposition, we compared these SOC data with previous litter measurements conducted using

110   comparable NMR methods (Supplementary Table 7). Although highly variable among

111   ecosystems, litter is typically dominated by carbohydrates (36–80%; mean 56%), with lesser

112   contributions from lignin (12–43%; mean 24%), lipids (0–21%; mean 8%) and proteins (0.3–

113   28%; mean 7%). Litter data available from a subset of the forested NEON sites also had a mean

114   lignin content of 24% (Extended Data Fig. 2). Therefore, on average, SOC from our surface

115   mineral soils (Fig. 1b) tended to have less carbohydrate, similar lignin, and more lipid and

116 protein as compared with typical organic matter inputs, even after accounting for carbohydrate

117 losses during HF treatment (Supplementary Information). Our data are consistent with the

118 emerging consensus that readily decomposable biomolecules, especially carbohydrates and

119 proteins, are often important SOC components[33], and they reinforce the importance of lipids[1,4].

120 However, our data challenge the view that lignin is disproportionately lost relative to other

121 molecules as microbes degrade litter to form SOC[5,14,15]. Although variable, mean lignin

122 abundance in SOC (21%) was similar to the other dominant molecules and similar to mean lignin

123 abundance in litter.

124       Some SOC molecules covaried with vegetation and management characteristics

125 (Extended Data Fig. 4). Lignin was significantly greater (23% vs. 16%) and protein was a

126 smaller component (15% vs. 23%) of SOC in forests than grasslands/shrublands ($P < 0.05$). Char

127 was greater in ecosystems experiencing periodic prescribed fire (22% vs. 16%, $P < 0.05$), but

128 intriguingly, char was not limited to fire-prone ecosystems. All soils contained measurable char

129 (> 6%). Interpretation of char-like C remains contentious, as it might be produced by non-fire-

130 related processes[34]. However, the fact that char significantly increased in soils with a known

131 history of fire (Extended Data Fig. 4) indicates the importance of pyrogenesis. Five soils were

132 from perhumid climates where mean annual precipitation (MAP) exceeded potential

133 evapotranspiration (PET) by > 1 m (Supplementary Table 1). This implicates a plausible role for

134 ancient or anthropogenic fire in producing extant char. For example, anthropogenic charcoal

135 production, but not natural fire, was documented in the rainforests of the Luquillo Mountains,

136 Puerto Rico[35], where two samples were collected.

137

138 **Consistent molecular tradeoffs linked to ecosystem factors**

139    Molecules covaried in predictable ways within samples despite high variability in composition

140    among samples. Lignin and carbonyl C were positively correlated while lignin and protein, and

141    carbohydrate and char, were each negatively correlated (corrected $P < 0.01$, $P < 0.0001$, and $P <$

142    $0.0001$, respectively; Extended Data Fig. 5). Principal components analysis showed that the

143    correlation matrix of SOC molecules was dominantly explained ($R^2 = 0.90$) by three axes, which

144    we rotated orthogonally to maximize interpretability and thus refer to as rotated components

145    (Fig. 2, Supplementary Table 3). The first rotated component (RC1) scores were positively

146    correlated with lignin ($r = 0.9$) and carbonyl ($r = 0.8$) and negatively correlated with protein ($r = $ -

147    $0.8$; Fig. 2a, Extended Data Fig. 6). The second rotated component (RC2) scores were positively

148    correlated with carbohydrate ($r = 0.88$) and negatively correlated with char ($r = -0.92$). The third

149    rotated component (RC3) scores were strongly correlated with lipid ($r = -0.98$) and weakly

150    correlated with other molecules ($r < 0.53$). The molecular tradeoffs implied by the RC axes

151    indicated the importance of multiple SOC persistence mechanisms enabling differential accrual

152    of molecules among ecosystems.

153         To assess potential mechanisms underlying observed variation in SOC composition, we

154    analyzed correlations between molecule relative abundance and biogeochemical predictors and

155    performed multiple regression analyses and structural equation models (SEMs) for each RC axis.

156    Several SOC molecules showed significant correlations with geochemical, biological, and

157    climate variables (Fig. 3, Extended Data Fig. 7, Extended Data Fig. 8). Lignin and carbonyl C

158    correlated positively with concentrations of oxalate-extractable aluminum ($Al_o$), which

159    represents Al in short-range-ordered (SRO) mineral phases and/or organo-metal complexes that

160    can protect SOC from microbial decomposition[36]. In contrast, protein had a negative correlation

161    with $Al_o$ and with copy numbers of the fungal internal transcribed spacer (ITS) region, and a

162 positive correlation with pH. Lipid C correlated negatively with mean annual temperature

163 (MAT), pH, and sulfate-extractable calcium and magnesium ($Ca_s + Mg_s$), which may participate

164 in cation bridging with clay minerals. Lipid C correlated positively with the thickness of the

165 overlying organic (O) horizon, which in turn had a strong negative relationship with MAT (Fig.

166 4).

167 Consistent with these pairwise correlations, different sets of variables best predicted

168 variation along each RC axis (Extended Data Fig. 9). We present multiple models fit by

169 backwards selection using more conservative ($P < 0.01$) and liberal ($P < 0.05$) variable selection

170 criteria, respectively (Methods). We also compared models fit to the NEON samples only vs. the

171 complete dataset, given that not all potential predictors were available for all samples (e.g., root

172 and microbial data). Across all models, $Al_o$ concentration was the best predictor of RC1 (r =

173 0.63–0.69, $P < 0.001$), with increasing values reflecting greater lignin vs. protein. The more

174 liberal models indicated that $Ca_s + Mg_s$, forest vegetation, and prescribed fire were also

175 positively correlated with RC1, as was ITS copy number. For RC2, crystalline iron mineral

176 concentration ($Fe_{d-o}$) was a consistently important predictor across models (r = 0.28–0.48, P <

177 0.01), which was associated with increased carbohydrate vs. char. The liberal models also

178 indicated a negative correlation of RC2 with mineral horizon thickness and fine root C:N, and a

179 positive correlation with fine root biomass. For RC3, soil pH and O horizon thickness were the

180 strongest predictors (r = 0.40–0.60, $P < 0.001$); MAP-PET and fine root C:N also correlated

181 positively with RC3. The more acidic soils with thicker O horizons were associated with greater

182 lipid relative abundance.

183 The SEMs showed that the strongest biogeochemical predictors of SOC composition

184 were ultimately related to climate, either directly, or via proxies for soil development which were

185  also related to climate (thickness of the O horizon and surface mineral genetic horizon; Fig. 4).

186  Concentrations of $Al_o$ increased with MAP-PET (excess moisture drives dissolution of Al-

187  bearing minerals[37]) and decreased with MAT. Similarly, $Fe_{d-o}$, which accumulates as soils

188  progressively weather, also increased with MAP-PET. Temperature impacted SOC composition

189  both directly and indirectly. Increasing MAT decreased O horizon thickness, consistent with

190  increased decomposition of unprotected organic matter with warmer temperature[38]. Organic

191  horizon thickness was directly linked to RC3, and also indirectly linked via effects on pH (O

192  horizons may promote acidification by leaching organic acids[39]). Thinner O horizons were also

193  associated with thinner mineral surface genetic horizons in our dataset, possibly reflecting

194  differences in soil profile development related to litter decomposition rates. Surface mineral

195  horizon thickness, in turn, was proximately linked to SOC composition (RC2). Including

196  vegetation type or fire did not improve any of the SEMs, possibly because these factors were

197  adequately reflected by climate or soil-horizon-related variables.

198      The relationships between SOC composition and biogeochemical predictors observed

199  here provide a molecular-level explanation for trends in SOC content among ecosystems noted

200  elsewhere. The concentration of $Al_o$ is among the best predictors of SOC content at local to

201  global scales[9,40], reflecting its formation of protective complexes with SOC[36]. Our data imply

202  that specific geochemical associations between lignin- and carbonyl-derived SOC and $Al_o$ could

203  explain increases in SOC content with $Al_o$ among soils. The finding that lignin was the only

204  molecule whose relative abundance significantly increased with SOC content (Fig. 3) also

205  accords with this interpretation. A strong $Al_o$-lignin relationship is consistent with previous

206  evidence of ligand exchange by carboxylated aromatics on SRO ordered Al phases[41] and high

207  concentrations of lignin-derived C observed in humid tropical soils rich in SRO phases[25,40,42].

208 While lignin was greater in forested than non-forested soils (Extended Data Fig. 4), the

209 relationship between lignin and $Al_o$—and ultimately, climate—was much stronger than the

210 relationship with vegetation (Figs. 3,4, Extended Data Table 3). Relationships between $Al_o$ and

211 lignin partially reconcile aspects of old and new SOM paradigms: lignin-derived C may

212 contribute significantly to SOC in some soils[12] (Fig. 1), but not because of inherent

213 recalcitrance[31,32]. Rather, lignin (and carbonyl C, which was strongly correlated with lignin) may

214 vary among ecosystems as a function of geochemical context. Our statistical models also

215 supported a role for Ca and Mg in protecting lignin; these cations can provide protective bridging

216 between anionic SOC functional groups and negatively charged mineral surfaces[9]. Minerals and

217 metals are effective predictors of SOC content[31,36], and interactions with specific SOC molecules

218 may underlie these patterns.

219

220 **Complementary mechanisms of SOC persistence**

221 The first observed tradeoff, between lignin and protein (Figs. 2,5), may reflect multiple

222 underlying mechanisms. First, where SRO mineral phases (i.e., $Al_o$) and physicochemical

223 protection are scarce, protein relative abundance may increase because it is a major microbial

224 biomass component that can be efficiently recycled between living and dead microbes, whereas

225 most lignin C is decomposed to carbon dioxide[18,21]. Second, in acidic soils, low abundance of

226 protein (Fig. 3) vs. lignin may be driven by inefficient litter decomposition and low microbial

227 necromass production[23]. Third, the negative correlation between fungal ITS copies and protein

228 (Fig. 3) suggests that fungi may play a role in the lignin-protein tradeoff. Fungi are dominant

229 decomposers of lignin[13] but have a higher biomass C:N (lower protein content) than bacteria[43].

230 Finally, the observed tradeoff between lignin and protein in SOC could reflect the fundamental

231    role of the lignin:N ratio in controlling litter decay rates. Protein is a dominant soil N pool, and

232    limited N availability to produce lignin-degrading enzymes could constrain lignin mass loss[44].

233          The second observed tradeoff, between carbohydrate and char (Figs. 2,4), was most

234    closely linked to $Fe_{d-o}$. Crystalline Fe phases are protective sorbents that may promote soil

235    aggregation[36,41,42], and these physicochemical protection mechanisms may explain increased

236    relative abundance of carbohydrate, an easily decomposed molecule. Increasing quality (lower

237    C:N ratios) and quantity of fine root biomass were also associated with greater carbohydrate. In

238    contrast, char became more abundant as $Fe_{d-o}$ and fine root quantity and quality decreased. We

239    interpret this second tradeoff as follows: where physicochemical protection is lacking, a complex

240    molecular structure involving a greater diversity of bond types and more stable bonds (such as

241    those contained in polyaromatic char-like SOC) becomes increasingly important. Accepting that

242    molecular structure alone cannot guarantee long-term persistence[32], the logistical challenges of

243    char decomposition[16,29] may increase its relative contribution to SOC where other protection

244    mechanisms are unavailable and root C inputs are small.

245          The third SOC axis was related most strongly to lipid content (Figs. 2,4). Temperature

246    and pH are known to impact SOC content of mineral soils[9,38], and our data indicate that this may

247    be influenced by lipid accrual (Figs. 3,4). Lipids were largely independent of other molecules

248    and increased in cold, acidic soils with thick overlying organic horizons, comprising up to 59%

249    of SOC. Constraints on microbial physiology may promote lipid persistence. Lipids are the most

250    chemically reduced constituents of SOC, requiring greater activation energy for oxidation than

251    other compounds[45]. Because decomposition reactions are temperature dependent, the Arrhenius

252    equation predicts that molecules with the highest activation energies (i.e. lipids) exhibit the

253    greatest increase in decomposition rate with increasing temperature if other protection

254 mechanisms are unavailable[38]. As such, accrual of lipids in cold mineral soils is consistent with

255 thermodynamic expectations.

256        Collectively, our continental-scale dataset supports a concise new framework for

257 understanding multiple complementary mechanisms of SOC persistence among ecosystems (Fig.

258 5). Debates as to the relative importance of microbial necromass vs. lignin in SOC[20,21,25,28] can be

259 reconciled in part by considering the biogeochemical heterogeneity of ecosystems: necromass

260 may be more important than lignin where reactive Al phases are scarce, and vice-versa.

261 Similarly, the contested role of molecular stability in SOC persistence[16,31,32] is also illuminated

262 by examination of broad ecosystem gradients: where physicochemical protection mechanisms

263 mediated by Fe are scarce and high-quality root C inputs are small, char assumes a more

264 important role. Finally, differences in temperature and pH among ecosystems were ultimately

265 linked to lipid abundance, informing debates as to which SOC forms may be most impacted by

266 near-term warming and acidification[38,46]. Over longer timescales, temperature and moisture

267 influence all three axes of this conceptual framework via soil development (Fig. 4). Collectively,

268 our data point to the power of a macrosystems approach in reconciling paradigmatic

269 controversies in SOC research.

270

271 **References**

272 1. Baldock, J. A., Masiello, C. A., Gélinas, Y. & Hedges, J. I. Cycling and composition of

273     organic matter in terrestrial and marine ecosystems. *Mar. Chem.* **92**, 39–64 (2004).

274 2. Sutton, R. & Sposito, G. Molecular structure in soil humic substances: the new view.

275     *Environ. Sci. Technol.* **39**, 9009–9015 (2005).

276   3.  Lehmann, J. & Kleber, M. The contentious nature of soil organic matter. *Nature* **528**, 60–68

277      (2015).

278   4.  Baldock, J. A. *et al.* Assessing the extent of decomposition of natural organic materials using

279      solid-state $^{13}$C NMR spectroscopy. *Aust. J. Soil Res.* **35**, 1061–1084 (1997).

280   5.  Mahieu, N., Randall, E. W. & Powlson, D. S. Statistical analysis of published carbon-13

281      CPMAS NMR spectra of soil organic matter. *Soil Sci. Soc. Am. J.* **63**, 307 (1999).

282   6.  Grandy, A. S. & Neff, J. C. Molecular C dynamics downstream: The biochemical

283      decomposition sequence and its impact on soil organic matter structure and function. *Sci.*

284      *Total Environ.* **404**, 297–307 (2008).

285   7.  Baldock, J. A. *et al.* Aspects of the chemical structure of soil organic materials as revealed

286      by solid-state $^{13}$C NMR spectroscopy. *Biogeochemistry* **16**, 1–42 (1992).

287   8.  Ahmad, R., Nelson, P. N. & Kookana, R. S. The molecular composition of soil organic

288      matter as determined by $^{13}$C NMR and elemental analyses and correlation with pesticide

289      sorption. *Eur. J. Soil Sci.* **57**, 883–893 (2006).

290   9.  Rasmussen, C. *et al.* Beyond clay: towards an improved set of variables for predicting soil

291      organic matter content. *Biogeochemistry* **137**, 297–306 (2018).

292   10. Cotrufo, M. F., Ranalli, M. G., Haddix, M. L., Six, J. & Lugato, E. Soil carbon storage

293      informed by particulate and mineral-associated organic matter. *Nat. Geosci.* **12**, 989–994

294      (2019).

295   11. Wagai, R. *et al.* Linking temperature sensitivity of soil organic matter decomposition to its

296      molecular structure, accessibility, and microbial physiology. *Glob. Change Biol.* **19**, 1114–

297      1125 (2013).

298    12. Waksman, S. A. & Iyer, K. R. N. Contribution to our knowledge of the chemical nature and

299        origin of humus: I. On the synthesis of the "humus nucleus". *Soil Sci.* **34**, 43–69 (1932).

300    13. Kirk, T. K. & Farrell, R. L. Enzymatic" combustion": the microbial degradation of lignin.

301        *Annu. Rev. Microbiol.* **41**, 465–501 (1987).

302    14. Amelung, W., Brodowski, S., Sandhage-Hofmann, A. & Bol, R. Combining biomarker with

303        stable isotope analyses for assessing the transformation and turnover of soil organic matter.

304        *Adv. Agron.* **100**, 155–250 (2008).

305    15. Thevenot, M., Dignac, M.-F. & Rumpel, C. Fate of lignins in soils: A review. *Soil Biol.*

306        *Biochem.* **42**, 1200–1211 (2010).

307    16. Bosatta, E. & Ågren, G. I. Soil organic matter quality interpreted thermodynamically. *Soil*

308        *Biol. Biochem.* **31**, 1889–1891 (1999).

309    17. Miltner, A., Bombach, P., Schmidt-Brücken, B. & Kästner, M. SOM genesis: microbial

310        biomass as a significant source. *Biogeochemistry* **111**, 41–55 (2011).

311    18. Cotrufo, M. F., Wallenstein, M. D., Boot, C. M., Denef, K. & Paul, E. The Microbial

312        Efficiency-Matrix Stabilization (MEMS) framework integrates plant litter decomposition

313        with soil organic matter stabilization: do labile plant inputs form stable soil organic matter?

314        *Glob. Change Biol.* **19**, 988–995 (2013).

315    19. Kallenbach, C. M., Frey, S. D. & Grandy, A. S. Direct evidence for microbial-derived soil

316        organic matter formation and its ecophysiological controls. *Nat. Commun.* **7**, 13630 (2016).

317    20. Ma, T. *et al.* Divergent accumulation of microbial necromass and plant lignin components in

318        grassland soils. *Nat. Commun.* **9**, 1–9 (2018).

319    21. Liang, C., Amelung, W., Lehmann, J. & Kästner, M. Quantitative assessment of microbial

320        necromass contribution to soil organic matter. *Glob. Change Biol.* **25**, 3578–3590 (2019).

321   22. Khan, K. S., Mack, R., Castillo, X., Kaiser, M. & Joergensen, R. G. Microbial biomass,

322       fungal and bacterial residues, and their relationships to the soil organic matter C/N/P/S

323       ratios. *Geoderma* **271**, 115–123 (2016).

324   23. Malik, A. A. *et al.* Land use driven change in soil pH affects microbial carbon cycling

325       processes. *Nat. Commun.* **9**, 3591 (2018).

326   24. Córdova, S. C. *et al.* Plant litter quality affects the accumulation rate, composition, and

327       stability of mineral-associated soil organic matter. *Soil Biol. Biochem.* **125**, 115–124 (2018).

328   25. Huang, W. *et al.* Enrichment of lignin-derived carbon in mineral-associated soil organic

329       matter. *Environ. Sci. Technol.* **53**, 7522–7531 (2019).

330   26. Wan, D. *et al.* Iron oxides selectively stabilize plant-derived polysaccharides and aliphatic

331       compounds in agricultural soils. *Eur. J. Soil Sci.* **70**, 1153–1163 (2019).

332   27. Hernes, P. J., Kaiser, K., Dyda, R. Y. & Cerli, C. Molecular trickery in soil organic matter:

333       hidden lignin. *Environ. Sci. Technol.* **47**, 9077–9085 (2013).

334   28. Klotzbücher, T., Kalbitz, K., Cerli, C., Hernes, P. J. & Kaiser, K. Gone or just out of sight?

335       The apparent disappearance of aromatic litter components in soils. *SOIL* **2**, 325–335 (2016).

336   29. Preston, C. M. & Schmidt, M. W. I. Black (pyrogenic) carbon: a synthesis of current

337       knowledge and uncertainties with special consideration of boreal regions. *Biogeosciences* **3**,

338       397–420 (2006).

339   30. Lehmann, J. *et al.* Australian climate–carbon cycle feedback reduced by soil black carbon.

340       *Nat. Geosci.* **1**, 832–835 (2008).

341   31. Mikutta, R., Kleber, M., Torn, M. S. & Jahn, R. Stabilization of soil organic matter:

342       association with minerals or chemical recalcitrance? *Biogeochemistry* **77**, 25–56 (2006).

343   32. Kleber, M. What is recalcitrant soil organic matter? *Environ. Chem.* **7**, 320–332 (2010).

344    33. Schmidt, M. W. I. *et al.* Persistence of soil organic matter as an ecosystem property. *Nature*

345      **478**, 49–56 (2011).

346    34. DiDonato, N., Chen, H., Waggoner, D. & Hatcher, P. G. Potential origin and formation for

347      molecular components of humic acids in soils. *Geochim. Cosmochim. Acta* **178**, 210–222

348      (2016).

349    35. Scatena, F. An introduction to the physiography and history of the Bisley experimental

350      watersheds in the Luquillo Mountains of Puerto Rico. *USDA For. Serv. Gen. Tech. Rep.* **SO-**

351      **72**, 1–22 (1989).

352    36. Kleber, M. *et al.* Mineral–organic associations: formation, properties, and relevance in soil

353      environments. *Adv. Agron.* **130**, 1–140 (2015).

354    37. Slessarev, E. W. *et al.* Water balance creates a threshold in soil pH at the global scale.

355      *Nature* **540**, 567–569 (2016).

356    38. Davidson, E. A. & Janssens, I. A. Temperature sensitivity of soil carbon decomposition and

357      feedbacks to climate change. *Nature* **440**, 165–173 (2006).

358    39. Lundström, U. S., van Breemen, N. & Bain, D. The podzolization process. A review.

359      *Geoderma* **94**, 91–107 (2000).

360    40. Kramer, M. G., Sanderman, J., Chadwick, O. A., Chorover, J. & Vitousek, P. M. Long-term

361      carbon storage through retention of dissolved aromatic acids by reactive particles in soil.

362      *Glob. Change Biol.* **18**, 2594–2605 (2012).

363    41. Kaiser, K. & Guggenberger, G. The role of DOM sorption to mineral surfaces in the

364      preservation of organic matter in soils. *Org. Geochem.* **31**, 711–725 (2000).

365    42. Coward, E. K., Ohno, T. & Plante, A. F. Adsorption and molecular fractionation of dissolved

366        organic matter on iron-bearing mineral matrices of varying crystallinity. *Environ. Sci.*

367        *Technol.* **52**, 1036–1044 (2018).

368    43. Throckmorton, H. M., Bird, J. A., Dane, L., Firestone, M. K. & Horwath, W. R. The source

369        of microbial C has little impact on soil organic matter stabilisation in forest ecosystems.

370        *Ecol. Lett.* **15**, 1257–1265 (2012).

371    44. Moorhead, D. L. & Sinsabaugh, R. L. A theoretical model of litter decay and microbial

372        interaction. *Ecol. Monogr.* **76**, 151–174 (2006).

373    45. LaRowe, D. E. & Van Cappellen, P. Degradation of natural organic matter: A

374        thermodynamic analysis. *Geochim. Cosmochim. Acta* **75**, 2030–2042 (2011).

375    46. Ye, C. *et al.* Reconciling multiple impacts of nitrogen enrichment on soil carbon: plant,

376        microbial and geochemical controls. *Ecol. Lett.* **21**, 1162–1173 (2018).

377

378    **Correspondence and requests for materials** should be addressed to S.J.H

379    (stevenjh@iastate.edu) and W.C.H. (william_hockaday@baylor.edu)

380

387

**Author contributions**

S.J.H., S.R.W., and W.C.H. developed the research concepts, C.Y. and W.C.H. conducted the

NMR analyses, S.J.H., C.Y., S.R.W, and W.C.H. analyzed data, and S.J.H. and C.Y. wrote the

paper with contributions from all authors.

**Competing interests**

The authors declare no competing interests.

**Additional Information**

**Supplementary Information** is available for this paper.

**Figure Captions**

**Fig. 1**: **Boxplots of carbon abundance as the fraction of total SOC in each sample.** Values

were determined directly from $^{13}$C CPMAS NMR peak areas (**a**) and by applying a molecular

mixing model (**b**). Grey dots represent observations (n = 42). Center lines are medians; box

limits are upper and lower quartiles; whiskers are 1.5x the interquartile ranges; points are

outliers.

**Fig. 2: Rotated principal components analysis of SOC molecules**. RC1, RC2, and RC3

represent rotated components 1–3, which respectively explained 35%, 29%, and 26% of the total

variation (90% overall) in the correlation matrix of SOC molecule relative abundance. Grey dots

represent soil samples, and labeled green arrows indicate correlations between SOC molecules

and RCs, with the correlation coefficient indicated on the top and right axes (carbohyd denotes

411 carbohydrate). Several samples with RC values > 3 (Supplementary Table 3) are not shown for

412 clarity.

413

414 **Fig. 3**: **Heatmap of correlations (r) between SOC molecules and biogeochemical predictors**.

415 The symbols *, **, and **** denote corrected significance at $P < 0.05$, $P < 0.01$, and $P < 0.0001$,

416 respectively. MAT, mean annual temperature; MAP-PET, mean annual precipitation minus

417 potential evapotranspiration. Additional descriptions of biogeochemical predictors and any data

418 transformations are provided in the Methods and Supplementary Table 5.

419

420 **Figure 4: Parsimonious structural equation models of SOC molecular composition.** The

421 response variables (RC1, RC2, RC3) are the rotated principal components of SOC molecule

422 relative abundance. Solid yellow and blue lines indicate significant positive or negative

423 piecewise relationships between variables ($P < 0.05$). Dashed lines indicate non-significant

424 piecewise relationships that improved overall model fit as indicated by comparing AIC of nested

425 models. Numbers in boxes are scaled correlation coefficients. Fisher's C statistic refers to the test

426 of the overall model fit, where high P values indicate plausibility of the overall model.

427

428 **Figure 5: Conceptual model of three-dimensional tradeoffs in SOC composition linked to**

429 **complementary persistence mechanisms as supported by our data.** Soil samples fall within a

430 spherical space indicating the relative predominance of different SOC molecules, which are

431 constrained according to three major axes of variation. The location of a sample along each axis

432 indicates the relative importance of different SOC persistence mechanisms as described in the

433 text.

434

**Methods**

**Soil sampling and analysis.** We analyzed the molecular SOC composition of surface mineral

soil samples spanning 32 sites in the NEON Megapit archive, along with 10 additional soils

which were selected to encompass additional diversity in biogeochemical characteristics

(Supplementary Table 1). Vegetation included forests (n = 29) and grasslands or open canopy

shrublands (n = 13), including both managed (burned or grazed) and wildland sites. Soils were

sampled from the upper-most mineral soil horizon at a given site (organic horizons were

excluded); complete soil profile descriptions for the NEON sites are provided in Supplementary

Table 2. Briefly, in the dominant soil and vegetation type at each NEON terrestrial site, a soil

profile was characterized and sampled by horizon with the help of US Department of Agriculture

Natural Resource Conservation (NRCS) staff and archived by NEON[47–49]. We requested

subsamples of A horizon material from each site in the Megapit archive that was available in

September 2019. The Gellisols had extensive organic (O) horizons, such that we requested

material from the mineral horizon closest to the surface (described as Bg/Oajj, A/Cjj, and Bg at

BONA, HEAL, and TOOL, respectively; Supplementary Table 1). The 10 non-NEON samples

analyzed here were each collected from 0–10 cm depth with a clean shovel after removing any

litter or O horizon material. All soils were air dried to constant mass and sieved to 2 mm. Visible

root fragments were removed with tweezers and soils were finely ground with a mortar and

pestle prior to subsequent analyses.

454

*¹³C CPMAS NMR analyses and sample preparation* All 42 samples were prepared for NMR

analyses, allowing a comparative characterization of organic C molecular composition. In order

457    to increase the NMR sensitivity and remove paramagnetic materials, soils were pre-treated with

458    hydrochloric acid (HCl, 10% wt.) and hydrofluoric acid (HF, 10%, wt.) to remove any calcium

459    carbonate and mineral phases, respectively[50]. Briefly, 2–3 g of finely ground soil was weighed

460    into a 50 mL sealed polyethylene centrifugation tube, saturated with 30 mL HCl, and allowed to

461    settle for 30 min. After centrifugation and discarding HCl, the remaining slurry was then shaken

462    with 40 ml of mixed HF (10% wt.) and HCl (10% wt.) for 8 h, and subsequently centrifuged. The

463    supernatant was removed and discarded appropriately. After repeating the procedure four times,

464    each sample was washed with distilled water three times and dried at 50 $^{\circ}$C under a stream of

465    dinitrogen gas.

466          Solid-state $^{13}$C CP-MAS and $^{13}$C DP-MAS NMR spectra were recorded at room

467    temperature (23 °C) using a 300 MHz Bruker AVANCE III NMR spectrometer equipped with a

468    4 mm magic angle spinning (MAS) probe (Bruker BioSpin, Billerica, MA) at Baylor University

469    (Waco, TX). The 60–130 mg HF-treated sample was placed in a zirconium rotor with a diameter

470    of 4 mm and Kel-F caps to maximize the C mass and signal intensity. A MAS rate of 12 kHz

471    was used for all NMR measurements. Cross polarization (CP) experiments used a ramped-

472    amplitude (50% to 100%) contact pulse and rotor synchronized Hahn echo[51]. The contact time

473    and recycle delay were set to 2 ms and 1.2 s, respectively, and composite pulse proton

474    decoupling was applied during signal acquisition. Direct polarization (DP) $^{13}$C spectra were

475    acquired with a 90-degree excitation pulse and rotor-synchronized Hahn echo[52], with a recycle

476    delay of 180 s. Glycine was used as an external standard for setting pulse angles, chemical shift

477    and Hartman-Hahn matching conditions. DPMAS spectra were obtained for 11 HF-treated soil

478    samples as a means against which to assess relative quantitation bias in CPMAS NMR data[53].

479    These samples were selected to span a broad range of biogeochemical diversity (nine soil orders;

480      Supplementary Table 1) and contained sufficient SOC (> 2.9% C in the original samples) to

481      enable timely analysis by DPMAS.

482         CPMAS spectra for HF-treated samples were acquired with more than 6000 scans. To

483      assess potential impacts of HF treatment on SOC composition, 11 untreated samples with

484      relatively high SOC concentration (> 6%) were also selected for NMR analysis (this set differed

485      slightly from the CPMAS/DPMAS comparison given the differing selection criteria). These

486      samples included six soil orders and spanned a broad range of paramagnetic element content

487      (14–58 mg Fe g$^{-1}$). Spectra for these untreated samples were recorded using the same operation

488      conditions of HF-treated samples and were acquired with more than 44000 scans. After baseline

489      correction, quantification was performed by dividing the spectra into seven chemical shift

490      regions: 0–45 ppm, 45–60 ppm, 60–95 ppm, 95–110 ppm, 110–145 ppm, 145–165 ppm and

491      165–215 ppm, assigned to alkyl C, N-alkyl + methoxyl C, O-alkyl C, Di-O-alkyl C, aromatic C,

492      phenolic C, amide + carbonyl C, respectively. Subsequently, a molecular mixing model was

493      applied to the seven integrated spectra regions, to estimate the relative abundances of six

494      molecular SOC constituents (carbohydrate, protein, lignin, lipid, carbonyl and char)[1]. The

495      elemental concentrations of C and N were measured on the HF-treated samples by

496      combustion/elemental analysis at Baylor University (Costech 4010, Valencia, CA) and were

497      used as additional constraints on the molecular mixing model solutions[1].

498

499      *Biogeochemical analyses* Megapit soil samples and vegetation in proximity to the soil pit were

500      subjected to numerous physical and chemical analyses[54]. Here, we utilized measurements of total

501      elemental content and particle size from the Megapit samples. We also used measurements of the

502      copy number of bacteria/archaea (16S) and fungi (ITS) coding regions calculated by quantitative

503     polymerase chain reaction (qPCR), which were conducted on separate fresh soil samples

504     collected in the vicinity of each sampled Megapit profile[54]. Briefly, these soils were flash frozen

505     in the field on dry ice and shipped to an analytical facility for DNA extraction and amplification.

506     Soil samples for qPCR analysis were collected periodically (approximately three times per year)

507     from each site from 0–30 cm depth, and cores were visually separated according to organic and

508     mineral horizons; only samples from mineral soil were used here. We selected samples from

509     plots in proximity to each Megapit (i.e., within several hundred m; denoted as Tower plots in

510     NEON terminology) and averaged the mean 16S and ITS abundance for each site based on the

511     2016–2018 data. Fine root biomass was measured by depth in three pit profiles within the

512     Megapit and sorted into live/dead classes for fine (< 2 mm or < 4 mm, depending on the site) and

513     coarse diameter classes. Here, we denoted the combined < 2 mm and < 4 mm fractions as fine

514     roots for subsequent analyses. Roots were dried, weighed, and combusted for analysis of carbon

515     (C) and nitrogen (N) content. We averaged root data from 0–30 cm depth for use in subsequent

516     analyses. No NEON root data were available from TOOL, so we used previous published data

517     from the same site[55]. Samples for foliar and/or litter chemistry were available from a subset of

518     the NEON sites (15 and 16 sites, respectively), as these are collected from each site on a five-

519     year rolling schedule. Foliar samples represented clips of bulk herbaceous samples from the plant

520     community. Litter samples included debris from trees and shrubs. We used measurements of

521     foliar and litter C:N and a proxy for lignin content (acid-unhydrolyzable residue)[54]. No root or

522     microbial or litter chemistry data were available from the non-NEON samples from which we

523     collected $^{13}C$ NMR spectra.

524         We conducted several additional soil extractions of all samples to quantify reactive

525     metals. Subsamples were extracted in parallel with sodium dithionite (1:150 ratio of

23

526    soil:solution) to quantify pedogenic iron (denoted $Fe_d$) and ammonium oxalate (1:60 ratio of

527    soil:solution) to quantify Fe and aluminum in short-range-ordered phases and organo-metal

528    complexes (termed $Fe_o$ and $Al_o$). The concentration of crystalline Fe minerals was then

529    calculated as the difference between $Fe_d$ and $Fe_o$ ($Fe_{d-o}$). Subsamples were also sequentially

530    extracted with deionized water and sodium sulfate (1:150 ratio of soil:solution). The calcium and

531    magnesium concentration of the sodium sulfate extraction (termed $Ca_s + Mg_s$), which followed

532    the water extraction, was interpreted as a proxy for Ca and Mg that may have participated in

533    divalent cation bridging between clays and organic matter[46]. All metals were analyzed by

534    inductively coupled plasma optical emission spectroscopy at Iowa State University (ICP-OES;

535    Perkin Elmer Optima 5300 DV, Waltham Massachusetts). Mean annual precipitation and

536    temperature data were estimated for each NEON site using previously synthesized data[56].

537    Potential evapotranspiration (PET) data were extracted from a global 1-km resolution mean

538    annual evapotranspiration dataset from 2000-2014[57].

539

540    *Statistical analyses* Correlation heatmaps were calculated between SOC molecules and

541    biogeochemical predictors, and some variables were log10-transformed because of skewness

542    ($Fe_o$, $Al_o$, $Fe_{d-o}$, $Ca_s + Mg_s$, ITS, 16S). Significance of correlations was calculated by multiplying

543    *P* values according to a Bonferroni correction to correct for multiple comparisons and an $\alpha$ of

544    0.10. We used a rotated principal components analysis to assess relationships among the relative

545    abundances of the six SOC molecules calculated from the molecular mixing model and soil

546    biogeochemical variables. Principal components were calculated from the correlation matrix of

547    the C molecule data and rotated orthogonally (varimax rotation) using the "Psych" package[58] in

548    R version 3.6.0. Rotation is commonly used in PCA to simplify interpretation of principal

549     components by maximizing/minimizing the correlations between factors and component axes.

550     These rotated components (RC) of the correlation matrix were used to facilitate interpretation of

551     each component in terms of dominant C molecule(s). To investigate relationships among RCs

552     and biogeochemical variables, we fit multiple linear regression models for each RC using the lm

553     function in R. The global model contained the following potential explanatory variables: mean

554     annual temperature, mean annual precipitation, mean annual precipitation minus potential

555     evapotranspiration, forest vs. non-forest vegetation, presence/absence of recurring fire, $Al_o$, $Fe_o$,

556     $Fe_{c-o}$, $Ca_s + Mg_s$, fine root biomass, fine root C:N ratio, the ratio of total base cations to

557     zirconium (a weathering ratio sensu[59]), and copy numbers of 16S and ITS genes quantified by

558     qPCR. Most of the non-forested ecosystems were grazed, such that a separate variable for

559     grazing was not included. Certain predictor variables were only available for the NEON samples

560     (n = 32; Supplementary Table 3), such that model selection was conducted independently for

561     both datasets. Candidate models for each RC were carefully investigated for multicollinearity of

562     predictors and assumptions of normality and heteroscedasticity by calculating variance inflation

563     factors (VIF) and graphically examining plots of residuals. Prior to model selection, individual

564     predictors with VIF > 3 were sequentially deleted[60], the reduced global model was refit, and VIF

565     values were calculated again. After removing collinear predictors, we performed model selection

566     by backwards elimination; more conservative and more liberal models yielded by $\alpha = 0.01$ and $\alpha$

567     $= 0.05$, respectively, were presented for completeness. Following multiple linear regression, to

568     better understand interrelationships among proximate predictors of RC axes and soil forming

569     factors, we fit SEMs using the piecewiseSEM package v. 2.1.0 in R[61]. Candidate SEMs included

570     the direct biogeochemical predictors identified by multiple linear regressions, along with climate,

571 soil, and vegetation variables that might influence those biogeochemical predictors. The

572 optimum models for each RC were selected by comparing AIC values among nested models.

573

574 **Data availability**

575 Summarized NMR data are available in the Supplementary Information, and raw NMR spectra

576 data and sample biogeochemical characteristics are available in the Environmental Data

577 Initiative digital repository: https://doi.org/10.6073/pasta/2284825ecb8460f056ae5b0e7d355cc8

578

579 **Code availability**

580 R scripts used for post-processing data are available in the Environmental Data Initiative digital

581 repository: https://doi.org/10.6073/pasta/2284825ecb8460f056ae5b0e7d355cc8

582

583 **References**

584 47. Ayres, E., *et al.* NEON Field and Lab Procedure and Protocol: TIS Soil Pit Sampling

585      Protocol. NEON.DOC.001307. https://data.neonscience.org/data-products/DP1.00097.001

586      (2017).

587 48. Ayres E., & Durden, D. NEON Field and Lab Procedure and Protocol: TIS Soil Archiving.

588      NEON.DOC.000325. https://data.neonscience.org/data-products/DP1.00097.001 (2017).

589 49. Ayres, E. NEON procedure and protocol: producing TIS soil archive subsamples for users.

590      NEON.DOC.001306. https://data.neonscience.org/data-products/DP1.00097.001 (2017).

591 50. Gélinas, Y., Baldock, J. A. & Hedges, J. I. Demineralization of marine and freshwater

592      sediments for CP/MAS $^{13}$C NMR analysis. *Org. Geochem.* **32**, 677–693 (2001).

593    51. Harbison, G. S. *et al.* High-resolution carbon-13 NMR of retinal derivatives in the solid state.

594        *J. Am. Chem. Soc.* **107**, 4809–4816 (1985).

595    52. Mao, J.-D. *et al.* Quantitative characterization of humic substances by solid-state carbon-13

596        nuclear magnetic resonance. *Soil Sci. Soc. Am. J.* **64**, 873–884 (2000).

597    53. Longbottom, T. L. & Hockaday, W. C. Molecular and isotopic composition of modern soils

598        derived from kerogen-rich bedrock and implications for the global C cycle. *Biogeochemistry*

599        **143**, 239–255 (2019).

600    54. NEON (National Ecological Observatory Network). DP1.00096.001, DP1.10066.001,

601        DP1.10102.001, DP1.10109.001 (accessed September 1, 2019), DP1.10026.001,

602        DP1.10033.001, DP1.10031.001 (accessed May 15, 2020). http://data.neonscience.org.

603    55. Sullivan, P. F. *et al.* Climate and species affect fine root production with long-term

604        fertilization in acidic tussock tundra near Toolik Lake, Alaska. *Oecologia* **153**, 643–652

605        (2007).

606    56. SanClements, M. *et al.* Collaborating with NEON. *BioScience* **70**, 107–107 (2020).

607    57. Mu, Q., Zhao, M. & Running, S. W. Improvements to a MODIS global terrestrial

608        evapotranspiration algorithm. *Remote Sens. Environ.* **115**, 1781–1800 (2011).

609    58. Revelle, W. psych: Procedures for Personality and Psychological Research, Northwestern

610        University, Evanston, Illinois, USA, https://CRAN.R-project.org/package=psych Version =

611        1.8.12 (2018).

612    59. Chittleborough, D. J. Indices of weathering for soils and palaeosols formed on silicate rocks.

613        *Aust. J. Earth Sci.* **38**, 115–120 (1991).

614    60. Hair, J. F., Risher, J. J., Sarstedt, M. & Ringle, C. M. When to use and how to report the

615        results of PLS-SEM. *Eur. Bus. Rev.* **31**, 2–24 (2019).

616    61. Lefcheck, J. S. piecewiseSEM: Piecewise structural equation modelling in R for ecology,

617        evolution, and systematics. *Methods Ecol. Evol.* **7**, 573–579 (2016).

618