

1 **Systematic variability enhances the reproducibility of an ecological study**

2 Alexandru Milcu^{1,2}, Ruben Puga-Freitas³, Aaron M. Ellison^{4,5}, Manuel Blouin^{3,6}, Stefan Scheu⁷,
3 Thomas Girin⁸, Grégoire T. Freschet², Laura Rose⁹, Michael Scherer-Lorenzen⁹, Sebastien
4 Barot⁶, Jean-Christophe Lata¹⁰, Simone Cesarz^{11,12}, Nico Eisenhauer^{11,12}, Agnès Gigon³,
5 Alexandra Weigelt^{11,12}, Amandine Hansart¹³, Anna Greiner⁹, Anne Pando⁶, Arthur Gessler^{14,15},
6 Carlo Grignani¹⁶, Davide Assandri¹⁶, Gerd Gleixner¹⁷, Jean-François Le Galliard^{10,13}, Katherine
7 Urban-Mead², Laura Zavattaro¹⁶, Marina E.H. Müller¹⁴, Markus Lange¹⁸, Martin Lukac^{19,20},
8 Michael Bonkowski¹⁷, Neringa Mannerheim²¹, Nina Buchmann²¹, Olaf Butenschoen^{7,22}, Paula
9 Rotter⁹, Rahme Seyhun¹⁹, Sebastien Devidal¹, Zachary Kayler^{14,23} and Jacques Roy¹

10 ¹Ecotron (UPS-3248), CNRS, Campus Baillarguet, F-34980, Montferrier-sur-Lez, France.

11 ²Centre d'Ecologie Fonctionnelle et Evolutive, CEFE-CNRS, UMR 5175, Université de
12 Montpellier – Université Paul Valéry – EPHE, 1919 route de Mende, F-34293, Montpellier
13 Cedex 5, France.

14 ³Institut des Sciences de l'Ecologie et de l'Environnement de Paris (UPMC, UPEC, Paris Diderot,
15 CNRS, IRD, INRA), Université Paris-Est Créteil, 61 avenue du Général De Gaulle, F-94010
16 Créteil Cedex, France.

17 ⁴Harvard Forest, Harvard University, 324 North Main Street, Petersham, Massachusetts, USA.

18 ⁵University of the Sunshine Coast, Tropical Forests and People Research Centre, Locked Bag 4,
19 Maroochydore DC, Queensland 4558, Australia.

20 ⁶IRD, Institut des Sciences de l'Ecologie et de l'Environnement de Paris (UPMC, UPEC, Paris
21 Diderot, CNRS, IRD, INRA), UPMC, Bâtiment 44-45, deuxième étage, bureau 208, CC 237, 4
22 place Jussieu, 75252 Paris cedex 05, France.

Milcu et al. 2016

23 ⁷J.F. Blumenbach Institute for Zoology and Anthropology, Georg August University Göttingen,
24 Berliner Str. 28, 37073 Göttingen, Germany.

25 ⁸Institut Jean-Pierre Bourgin, INRA, AgroParisTech, CNRS, Université Paris-Saclay, RD10,
26 78026 Versailles Cedex, France.

27 ⁹Faculty of Biology, University of Freiburg, Geobotany, Schaenzlestr. 1, D-79104 Freiburg,
28 Germany.

29 ¹⁰Institut des Sciences de l'Ecologie et de l'Environnement de Paris (UPMC, UPEC, Paris
30 Diderot, CNRS, IRD, INRA), Sorbonne Universités, CC 237, 4 place Jussieu, 75252 Paris cedex
31 05, France.

32 ¹¹German Centre for Integrative Biodiversity Research (iDiv), Halle-Jena-Leipzig, Deutscher
33 Platz 5e, 04103 Leipzig, Germany.

34 ¹²Institute of Biology, Leipzig University, Johannisallee 21, 04103 Leipzig, Germany.

35 ¹³Ecole normale supérieure, PSL Research University, Département de biologie, CNRS, UMS
36 3194, Centre de recherche en écologie expérimentale et prédictive (CEREPEP-Ecotron
37 IleDeFrance), 78 rue du château, 77140 Saint-Pierre-lès-Nemours, France.

38 ¹⁴Leibniz Centre for Agricultural Landscape Research (ZALF), Institute of Landscape
39 Biogeochemistry, Eberswalder Str. 84, 15374 Müncheberg, Germany.

40 ¹⁵Swiss Federal Research Institute WSL, Zürcherstr. 111, 8903 Birmensdorf, Switzerland.

41 ¹⁶Department of Agricultural, Forest and Food Sciences, University of Turin, largo Braccini, 2,
42 10095 Grugliasco, Italy.

43 ¹⁷Department of Terrestrial Ecology, Institute for Zoology, University of Cologne, Zulpicher Str.
44 47b, 50674 Köln, Germany.

45 ¹⁸Max Planck Institute for Biogeochemistry, Postfach 100164, 07701 Jena, Germany.

Milcu et al. 2016

46 ¹⁹School of Agriculture, Policy and Development, University of Reading, Reading, RG6 6AR,
47 UK.

48 ²⁰FLD, Czech University of Life Sciences, 165 00 Prague, Czech Republic.

49 ²¹Institute of Agricultural Sciences, ETH Zurich, Universitätsstrasse 2, 8092 Zürich, Switzerland

50 ²²Senckenberg Biodiversität und Klima Forschungszentrum BiK-F, Georg-Voigt-Straße 14-16,
51 Frankfurt am Main.

52 ²³USDA Forest Service, Northern Research Station, Lawrence Livermore National Laboratory,
53 Livermore, California 94550 USA.

54

55 **Corresponding author:** Alexandru Milcu, CNRS, Ecotron - UPS 3248, Campus Baillarguet, 34980,
56 Montferrier-sur-Lez, France, email: alex.milcu@cnr.fr, phone: +33 (0) 434-359-893.

57

58 **Many scientific disciplines currently are experiencing a “reproducibility crisis” because**
59 **numerous scientific findings cannot be repeated consistently¹⁻⁴. A new but controversial**
60 **hypothesis postulates that stringent levels of environmental and biotic standardization in**
61 **experimental studies reduces reproducibility by amplifying impacts of lab-specific**
62 **environmental factors not accounted for in study designs⁵⁻⁸. A corollary to this hypothesis**
63 **is that the deliberate introduction of controlled systematic variability (CSV) in**
64 **experimental designs can increase reproducibility. We tested this hypothesis using a multi-**
65 **laboratory microcosm study in which the same ecological experiment was repeated in 14**
66 **laboratories. Each laboratory introduced environmental and genotypic CSV within and**
67 **among treatments in replicated microcosms established in either growth chambers (with**
68 **stringent control of environmental conditions) or glasshouses (with more variable**

69 **environmental conditions). The introduction of genotypic CSV increased reproducibility of**
70 **results in growth chambers but had no significant effect in glasshouses where**
71 **reproducibility also was lower. Environmental CSV had little effect on reproducibility.**
72 **This first deliberate attempt at reproducing an ecological experiment with added CSV**
73 **reveals that introducing genotypic CSV in experiments carried out under controlled**
74 **environmental conditions with stringent standardization can increase reproducibility by**
75 **buffering against unaccounted lab-specific environmental and biotic factors that may**
76 **otherwise strongly bias experimental outcomes.**

77

78 **Keywords:** standardization, microcosms, “reproducibility crisis”, experimental methods,
79 genotypic diversity, controlled environment, growth chambers, greenhouses, systematic
80 heterogenization

81

82 Reproducibility—the ability to duplicate a study and its findings—is a defining feature of
83 scientific research. In ecology, it is often argued that it is virtually impossible to precisely
84 duplicate any single ecological experiment or observational study because complex ecological
85 interactions between the ever-changing environment and the extraordinary diversity of biological
86 systems exhibiting a wide range of plastic responses at different levels of biological organization
87 together make exact duplication unfeasible^{9,10}. Although this may be true for observational and
88 field studies, numerous ecological (and agronomic) studies are carried out with artificially
89 assembled, simplified ecosystems and controlled environmental conditions in experimental
90 microcosms or mesocosms (henceforth, “microcosms”)^{11–13}. Since biotic and environmental
91 parameters can be tightly controlled in microcosms, results from such studies should be easier to

92 reproduce. Even though microcosms frequently have been used to address fundamental
93 ecological questions^{12,14,15}, there has been no quantitative assessment of the reproducibility of
94 any microcosm experiment.

95 Because it reduces within-treatment variability, experimental standardization—the
96 implementation of strictly defined and controlled properties of organisms and their
97 environment—is widely thought to increase both reproducibility and the sensitivity of statistical
98 tests^{7,16}. This paradigm has been challenged recently by several studies on animal behavior that
99 suggest that stringent standardization may, counterintuitively, be responsible for generating non-
100 reproducible results⁵⁻⁷; the results may be valid under given conditions (i.e., they are local
101 “truths”) but are not generalizable^{16,17}. Despite rigorous adherence to experimental protocols,
102 laboratories inherently vary in many conditions that are not measured and are thus unaccounted
103 for, such as experimenter, micro-scale environmental heterogeneity, physico-chemical properties
104 of reagents and lab-ware, pre-experimental conditioning of organisms, and their genetic and
105 epigenetic variation. It even has been suggested⁵⁻⁷ that attempts to stringently control all sources
106 of biological and environmental variation might inadvertently lead to the amplification of these
107 unmeasured variations among laboratories, thus reducing reproducibility. Some studies have
108 gone even further, hypothesizing that the introduction of controlled systematic variation (CSV)
109 among the replicates of a treatment (e.g., using different genotypes for different experimental
110 replicates or varying pre-experimental conditions) should lead to less variable mean response
111 values between the laboratories that duplicated the experiments^{6,7}. In short, reproducibility
112 should increase by shifting the variance from among experiments to within them⁷. If this is true,
113 then introducing CSV will increase researchers’ abilities to draw generalizable conclusions about

114 the directions and effect sizes of experimental treatments, while at the same time reducing the
115 probability of detecting statistically significant treatment effects.

116 To test the hypothesis that introducing CSV enhances reproducibility in an ecological
117 context, we had 14 European laboratories simultaneously run a simple microcosm experiment
118 using grass (*Brachypodium distachion* L.) monocultures and grass and legume (*Medicago*
119 *truncatula* Gaertn.) mixtures. This experiment measured the effects of the presence of a nitrogen-
120 fixing legume on ecosystem functioning and productivity in grass-legume mixtures ('net legume
121 effect' hereafter), an approach often used in legume-grass binary cropping systems^{18,19} and
122 biodiversity-ecosystem function experiments^{20,21}. All laboratories were provided with the same
123 experimental protocol, seed stock from the same batch, and identical containers with which to
124 establish microcosms with grass only and grass-legume mixtures. Alongside a control (CTR)
125 with no CSV and containing a homogenized soil substrate (mixture of soil and sand) and a single
126 genotype of each plant species, we explored the effects of five different types of within- and
127 among-microcosm CSV on experimental reproducibility of the net legume effect (Fig.1): 1)
128 within-microcosm environmental CSV (CSV-WE) achieved by spatially varying soil resource
129 distribution through the introduction of six sand patches into the soil; 2) among-microcosm
130 environmental CSV (CSV-AE), which varied the number of sand patches (none, three or six)
131 among replicate microcosms; 3) within-microcosm biological CSV (CSV-WB) that used three
132 distinct genotypes per species planted in homogenized soil in each microcosm; 4) among-
133 microcosm biological CSV (CSV-AB) that varied the number of genotypes (one, two or three)
134 planted in homogenized soil among replicate microcosms; and 6) both environmental and biotic
135 CSV (CSV-WEB) within microcosms that used six sand patches and three plant genotypes per
136 species in each microcosm. In addition, we tested whether CSV effects depended on the level of

137 standardization within laboratories by using two common experimental approaches ('SETUP'
138 hereafter): growth chambers with tightly controlled environmental conditions and identical soil
139 (eight laboratories) or glasshouses with more loosely controlled environmental conditions and
140 different soils (six laboratories; Extended Data Table 1). We first tested the response to CSV of
141 twelve variables that are used commonly to describe ecosystem functions of plant-soil
142 microcosms (Extended Data Table 2). We then determined how the different types of CSV
143 affected the mean effect size and its standard deviation (SD) within and among laboratories;
144 lower among-laboratory SD implies that the results were reproduced more closely.

145 Although each laboratory followed the same experimental protocol, we found remarkably
146 high levels of among-laboratory variation in mean values for the majority of response variables
147 and the net legume effect on those variables (Extended Data Figs 1 and 2). For example, the net
148 legume effect on mean total plant biomass varied from 1.31 to 6.72 g dry weight (DW) per
149 microcosm among growth chambers, suggesting that unmeasured laboratory-specific conditions
150 outweighed effects of experimental standardization. Among glasshouses, differences were even
151 larger: mean plant biomass varied by nearly two orders of magnitude, from 0.14 to 14.57g DW
152 per microcosm (Extended Data Fig. 2).

153 Among-laboratory SD of net legume effect was significantly affected by CSV, SETUP and
154 their interaction (Table 1, Fig. 2a, b and Extended Data Fig. 3). The main effect of CSV was the
155 lower values of among-laboratory SD in the CSV-AB treatment level relative to CTR ($t_{1,45} =$
156 1.97, $P = 0.054$, Extended Data Fig. 4b), indicating increased reproducibility for CSV-AB. The
157 CSV \times SETUP interaction was reflected in the result that among-laboratory SD for CSV-AB was
158 significantly lower only for growth chambers ($t_{1,21} = 2.40$, $P = 0.025$, Fig. 2a) and not for
159 glasshouses. Assuming that the grand mean (mean of all laboratories and CSV treatment levels)

160 is the best available estimate of the “true” legume effect, we also assessed how the CSV
161 treatment affected the deviation from the grand mean (Extended Data Fig. 5). We found that of
162 the five types of CSV, CSV-AB (among-microcosm variance in genotypes) differed least from
163 the grand mean and resulted in the most reproducible results (Fig. 2c).

164 Within-laboratory SD of the net legume effect was only marginally affected by CSV
165 treatment when the analysis was performed on within-laboratory SD from individual variables
166 (Table 1), but this effect was significant when the analysis was performed on the second
167 principal component (PC2) of a PCA analysis that included all twelve response variables
168 (Extended Data Table 3). No significant CSV \times SETUP interaction was found (Fig. 3a).
169 However, we did observe a significant SETUP effect (Table 1 and Extended Data Table 3):
170 within-laboratory SD was lower in growth chambers (Fig. 3b). As we observed a tendency for
171 CSV to increase within-laboratory variation, we also analyzed the impact of the most
172 reproducible CSV treatment—CSV-AB—on the statistical power of detecting the net legume
173 effect within individual laboratories. Adding CSV-AB led to a reduction in statistical power
174 (57% in CTR vs. 45% in CSV-AB) that could be compensated for by doubling the number of
175 microcosms per treatment.

176 We further explored the relationship between within- and among-laboratory SD to
177 determine whether reproducibility was increased by shifting the variation from among to within
178 laboratories. Although the introduction of CSV generally increased within-laboratory SD of the
179 net legume effect (Extended Data Fig. 6), the treatment level with the highest reproducibility
180 (CSV-AB in growth chambers) only exhibited a non-significant trend of higher within-laboratory
181 SD relative to CTR (Fig. 3c). Moreover, a statistical model of among-laboratory SD as a function
182 of within-laboratory SD, SETUP, and CSV treatment did not reveal a significant three-way

183 (within-laboratory SD \times SETUP \times CSV) interaction ($F_{5,120} = 0.49$, $P = 0.784$), although in the
184 growth chamber setup the steepest and flattest slopes were for the CTR and CSV-AB treatment
185 levels, respectively (Fig. 3c). Therefore, although the observed trends are in line with the
186 proposed conjecture that adding CSV enhances reproducibility by increasing within-laboratory
187 variability, our results do not provide unequivocally support for it.

188 Overall, our findings provide compelling support for the hypothesis that introducing CSV in
189 experimental designs can increase reproducibility of ecological studies⁵⁻⁷. We also suggest that
190 the relationship between CSV and reproducibility is purely probabilistic and results from the
191 decreased likelihood that microcosms containing CSV will respond to unaccounted lab-specific
192 environmental factors in the same direction and with the same magnitude. In particular,
193 introducing CSV by using multiple genotypes of study species among replicated microcosms
194 appears to be a good strategy to enhance reproducibility because a mixture of genotypes may
195 limit genotype-specific or assemblage-wide responses to lab-specific environmental variation
196 (see also additional discussion in Supplementary Information). This suggestion is in line with
197 mounting evidence that ecological responses differ significantly among genotypes, and that
198 failure to account for genetic diversity leads to spurious results²²⁻²⁴. Interestingly, the
199 effectiveness of CSV-AB in increasing reproducibility was higher in growth chambers, and this
200 result amplifies the importance of introducing CSV in designs with stringent environmental
201 standardization. The lack of a significant effect on reproducibility when introducing CSV-AB in
202 glasshouses could be explained by spatially and temporarily more heterogeneous conditions in
203 the glasshouse environment (e.g. light fleck effects, temperature gradients, etc.) that likely added
204 additional random/unsystematic variation within laboratories. We conclude that to increase
205 reproducibility, ecological experiments should include both rigorous standardization and

206 controlled systematic variability. Although there are multiple causes for the “reproducibility
207 crisis”^{4,25,26}, here we show that deliberately including genetic variation may be the simplest
208 solution for increasing the reproducibility of ecological studies performed in controlled
209 environments.

210

211 **References**

- 212 1. Massonnet, C. *et al.* Probing the reproducibility of leaf growth and molecular phenotypes:
213 a comparison of three Arabidopsis accessions cultivated in ten laboratories. *Plant Physiol.*
214 **152**, 2142–2157 (2010).
- 215 2. Begley, C. G. & Ellis, M. L. Raise standards for preclinical cancer research. *Nature* **483**,
216 531–533 (2012).
- 217 3. Open Science Collaboration. Estimating the reproducibility of psychological science.
218 *Science* **349**, aac4716 (2015).
- 219 4. Parker, T. H. *et al.* Transparency in ecology and evolution: real problems, real solutions.
220 *Trends Ecol. Evol.* **31**, 711–719 (2016).
- 221 5. Richter, S. H., Garner, J. P. & Würbel, H. Environmental standardization: cure or cause of
222 poor reproducibility in animal experiments? *Nat. Methods* **6**, 257–261 (2009).
- 223 6. Richter, S. H., Garner, J. P., Auer, C., Kunert, J. & Würbel, H. Systematic variation
224 improves reproducibility of animal experiments. *Nat. Methods* **7**, 167–8 (2010).
- 225 7. Richter, S. H. *et al.* Effect of population heterogenization on the reproducibility of mouse
226 behavior: a multi-laboratory study. *PLoS One* **6**, e16461 (2011).
- 227 8. Wolfinger, R. D. Reanalysis of Richter et al. (2010) on reproducibility. *Nat. Methods* **10**,

- 228 373–4 (2013).
- 229 9. Cassey, P. & Blackburn, T. Reproducibility and Repeatability in Ecology. *Bioscience* **56**,
230 958–9 (2006).
- 231 10. Ellison, A. M. Repeatability and transparency in ecological research. *Ecology* **91**, 2536–
232 2539 (2010).
- 233 11. Lawton, J. H. The Ecotron facility at Silwood Park: the value of ‘big bottle’ experiments.
234 *Ecology* **77**, 665–669 (1996).
- 235 12. Benton, T. G., Solan, M., Travis, J. M. & Sait, S. M. Microcosm experiments can inform
236 global ecological problems. *Trends Ecol. Evol.* **22**, 516–521 (2007).
- 237 13. Drake, J. M. & Kramer, A. M. Mechanistic analogy: how microcosms explain nature.
238 *Theor. Ecol.* **5**, 433–444 (2012).
- 239 14. Fraser, L. H. & Keddy, P. The role of experimental microcosms in ecological research.
240 *Trends Ecol. Evol.* **12**, 478–481 (1997).
- 241 15. Srivastava, D. S. *et al.* Are natural microcosms useful model systems for ecology? *Trends*
242 *Ecol. Evol.* **19**, 379–384 (2004).
- 243 16. De Boeck, H. J. *et al.* Global change experiments: challenges and opportunities.
244 *Bioscience* (2015). doi:10.1093/biosci/biv099
- 245 17. Moore, R. P. & Robinson, W. D. Artificial bird nests, external validity, and bias in
246 ecological field studies. *Ecology* **85**, 1562–1567 (2004).
- 247 18. Sleugh, B., Moore, K. J., George, J. R. & Brummer, E. C. Binary Legume–Grass Mixtures
248 Improve Forage Yield, Quality, and Seasonal Distribution. *Agron. J.* **92**, 24–29 (2000).
- 249 19. Nyfeler, D., Huguenin-Elie, O., Suter, M., Frossard, E. & Lüscher, A. Grass-legume
250 mixtures can yield more nitrogen than legume pure stands due to mutual stimulation of

- 251 nitrogen uptake from symbiotic and non-symbiotic sources. *Agric. Ecosyst. Environ.* **140**,
252 155–163 (2011).
- 253 20. Temperton, V. M., Mwangi, P. N., Scherer-Lorenzen, M., Schmid, B. & Buchmann, N.
254 Positive interactions between nitrogen-fixing legumes and four different neighbouring
255 species in a biodiversity experiment. *Oecologia* **151**, 190–205 (2007).
- 256 21. Suter, M. *et al.* Nitrogen yield advantage from grass-legume mixtures is robust over a
257 wide range of legume proportions and environmental conditions. *Glob. Chang. Biol.* **21**,
258 2424–2438 (2015).
- 259 22. Reusch, T. B., Ehlers, A., Hämmerli, A. & Worm, B. Ecosystem recovery after climatic
260 extremes enhanced by genotypic diversity. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 2826
261 (2005).
- 262 23. Noguera, D. *et al.* Amplifying the benefits of agroecology by using the right cultivars.
263 *Ecol. Appl.* **23**, 515–522 (2013).
- 264 24. Hughes, A. R., Inouye, B. D., Johnson, M. T., Underwood, N. & Vellend, M. Ecological
265 consequences of genetic diversity. *Ecol. Lett.* **11**, 609–623 (2008).
- 266 25. Nuzzo, R. How scientists fool themselves – and how they can stop. *Nature* **526**, 182–185
267 (2015).
- 268 26. Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* **533**, 452–454 (2016).
- 269

270 **Acknowledgements**

271 This study benefited from the CNRS human and technical resources allocated to the
272 ECOTRONS Research Infrastructures and the state allocation 'Investissement d'Avenir' ANR-
273 11-INBS-0001 as well as from financial support by the ExpeER (grant no. 262060) consortium

274 funded under the EU-FP7 research program (FP2007-2013). *Brachypodium* seeds were kindly
275 provided by Richard Sibout (Observatoire du Végétal, Institut Jean-Pierre Bourgin, F-78026
276 Versailles Cedex France) and *Medicago* seeds were supplied by Jean-Marie Prosperi (INRA
277 Biological Resource Centre, F-34060 Montpellier Cedex 1 France). We further thank Jean
278 Varale, Gesa Hoffmann, Paul Werthenbach, Oliver Ravel, Clement Piel and Damien Landais for
279 assistance throughout the study. For additional acknowledgements see Supplementary
280 Information.

281 **Author contributions**

282 A.M. and J.R. designed the study with input from M.B, S.B and J-C.L. Substantial methodological
283 contributions were provided by M.B., S.S, T.G., L.R. and M.S-L. Conceptual feedback on an early
284 version was provided by G.F., N.E., J.R. and A.M.E. Data were analysed by A.M. with input from
285 A.M.E. A.M. wrote the manuscript with input from all co-authors. All co-authors were involved
286 in carrying out the experiments and/or analyses.

287 **Author Information**

288 The authors declare no conflict of interest. Correspondence and request for materials should be
289 addressed to A.M. (alex.milcu@cnr.fr).

290 **METHODS**

291 All laboratories tried to the best of their abilities to carry out an identical experimental protocol.
292 Whereas not all laboratories managed to recreate precisely all details of the experimental
293 protocol, we considered this to be a realistic scenario under which ecological experiments using
294 microcosms are performed in glasshouses and growth chambers.

295 **Germination.** The seeds from the three genotypes of *Brachypodium distachyon* (Bd21, Bd21-3
296 and Bd3-1) and *Medicago truncatula* (L000738, L000530 and L000174) were first sterilized by
297 soaking 100 seeds in 100 mL of sodium hypochlorite solution at 2.6% of active chlorine and
298 stirred for 15 min using a magnet. Thereafter, the seeds were rinsed 3 times in 250mL of sterile
299 water for 10-20 seconds under shaking. Sterilized seeds were germinated in trays (10 cm depth)
300 filled with vermiculite. The trays were first kept at 4°C in the dark for three days before they
301 were moved to light conditions (300 $\mu\text{mol m}^{-2} \text{s}^{-1}$ PAR) and 20/16°C and 60/70% air RH for day-
302 and night-time, respectively. When the seedlings of both species reached 1 cm in height above
303 the vermiculite they were transplanted into the microcosms.

304 **Preparation of microcosms**

305 All laboratories used identical containers (2-liter volume, 14.8-cm diameter, 17.4-cm height).
306 Sand patches were created using custom-made identical “patch makers” consisting of six rigid
307 PVC tubes of 2.5-cm diam. and 25-cm length, arranged in a circular pattern with an outer
308 diameter of 10cm. A textile mesh was placed at the bottom of the containers to prevent the
309 spilling of soil through drainage holes. Filling of microcosms containing sand patches started
310 with the insertion of the “patch maker” into containers. Thereafter, in growth chamber setups,
311 2000 g dry weight of soil, subtracting the weight of the sand patches, was added into the
312 containers and around the tubes of the “patch maker”. In the glasshouse setups with different
313 soils, the dry weight of the soil differed slightly (depending on the soil density) and was first
314 estimated individually in each laboratory as the amount of soil we needed to fill the pots up to 2
315 cm from the top. Finally, the tubes were filled with a mixture of 10% soil and 90% sand. When
316 the microcosms did not contain sand patches, the amount of sand contained in six patches was
317 homogenized with the soil. During the filling of the microcosms, a common substrate for

318 measuring litter decomposition was inserted at the center of the microcosm at 8-cm depth. For
319 simplicity as well as for its fast decomposition rate, we used a single batch of commercially
320 available tetrahedron-shaped synthetic tea bags (mesh size of 0.25 mm) containing 2 g of green
321 tea (Lipton, Unilever), as proposed by the “tea-bag index” method²⁷. Once filled, the microcosms
322 where watered until water could be seen pouring out of the pot. The seedlings were then
323 manually transplanted to predetermined positions (Fig. 1), depending on the genotype and
324 treatment. Each laboratory established two blocks of 36 microcosms each, resulting in a total of
325 72 microcosms per laboratory, with blocks representing two distinct chambers in growth
326 chamber setups or two distinct growth benches in the same glasshouse.

327 **Soils**

328 All laboratories using growth chamber setups used the same soil, whereas the laboratories using
329 glasshouses used different soils (see Extended Data Table 1 for the physicochemical properties
330 of the soils). The soil used in growth chambers was classified as a nutrient-poor cambisol and
331 was collected from the top layer (0–20 cm) of a natural meadow at the Centre de Recherche en
332 Ecologie Expérimentale et Prédictive—CEREEP (Saint-Pierre-Lès-Nemours, France). Soils used
333 in glasshouses originated from different locations. The soil used by laboratory L2 was a fluvisol
334 collected from the top layer (0-40 cm) of a quarry site near Avignon, in the Rhône valley,
335 Southern France. The soil used by laboratory L4 was collected from near the La Cage field
336 experimental system (Versailles, France) and was classified as a luvisol. The soil used in the
337 glasshouse experiments L11 and L12 was collected from the top layer (0-20cm) within the haugh
338 of the river Dreisam in the East of Freiburg, Germany. This soil was classified as an umbric
339 gleysol with high organic carbon content. The soil from laboratory L14 was classified as a eutric
340 fluvisol and was collected on the field site of the Jena Experiment, Germany. Prior to the

341 established of microcosms, all soils were air dried at room temperature for several weeks and
342 sieved with a 2-mm mesh sieve. A common inoculum was provided to all laboratories to assure
343 that rhizobia specific to *M. truncatula* were present in all soils.

344 **Abiotic environmental conditions**

345 The set points for environmental conditions were 16-hour light (at $300 \mu\text{mol m}^{-2} \text{s}^{-1}$ PAR) and 8-
346 hour dark, 20/16 °C, 60/70% air RH for day- and night-time, respectively. Different soils (for
347 glasshouses) and treatments with sand patches likely affected water drainage and
348 evapotranspiration. The watering protocol was thus based on drying weight relative to weight at
349 full water holding capacity (WHC). The WHC was estimated based on the weight difference
350 between the dry weight of the containers and the wet weight of the containers 24 h after
351 abundant watering (until water was flowing out of the drainage holes in the bottom of each
352 container). Soil moisture was maintained between 60 and 80% of WHC (i.e. the containers were
353 watered when the soil water dropped below 60% of WHC and water added to reach 80% of
354 WHC) during the first 3 weeks after seedling transplantation and between 50 and 70% of WHC
355 for the rest of the experiment. Microcosms were watered twice a week with estimated WHC
356 values from two microcosms per treatment. To ensure that that the patch/heterogeneity
357 treatments did not become a water availability treatment, all containers were weighed and
358 brought to 70 or 80% of WHC every two weeks. This operation was synchronized with within-
359 block randomization. All 14 experiments were performed between October 2014 and March
360 2015.

361 **Sampling and analytical procedures**

362 After 80 days, the experiments were stopped and all plants were harvested. Plant shoots were cut
363 at the soil surface level, separated into species and dried at 60°C for three days. Roots and the

364 remaining litter in the tea bags were washed out of the soil using a 1-mm mesh sieve and dried at
365 60°C for three days. Microcosm evapotranspiration rate was measured before the harvesting as
366 the difference in weight changes from 70% of WHC after 48h. Shoot %C, %N, $\delta^{13}\text{C}$, and $\delta^{15}\text{N}$
367 were measured on pooled shoot biomass (including seeds) of *B. distachyon* and analyzed at the
368 Göttingen Centre for Isotope Research and Analysis using a coupled system consisting of an
369 elemental analyzer (NA 1500, Carlo Erba, Milan, Italy) and a gas isotope mass spectrometer
370 (MAT 251, Finnigan, Thermo Electron Corporation, Waltham, Massachusetts, USA).

371 **Data analysis and statistics**

372 We focused our analyses on the net legume effect—the difference between the equivalent
373 microcosms with and without legumes—as we considered that comparing within- and among-
374 laboratory variation in the effect size of an experimental treatment (here the presence of a
375 legume) was a more realistic test of reproducibility than comparing absolute values of response
376 variables. All analyses were performed using R version 3.2.4²⁸. To assess reproducibility we
377 investigated how CSV treatments affected the standard deviation (SD) of the measured variables,
378 with lower among-laboratory SD indicating increased reproducibility. We opted for SD instead
379 of the coefficient of variation because the net legume effect contained both positive and negative
380 values. As a complementary approach to assess the impact of CSV on reproducibility we
381 explored the extent to which the net legume effect was different from the grand mean (pooled
382 across all laboratories, CSV treatments, and two SETUPS) and used a Kruskal-Wallis test on the
383 ranked differences (of all response variables) from the grand mean.

384 Among-laboratory SD was computed from laboratory means for each response variable,
385 CSV treatments and SETUPS ($n = 144$; 6 CSV levels \times 2 SETUP levels \times 12 response variables).
386 Some of the twelve response variables are intrinsically correlated, but most did not have

387 correlation coefficients > 0.5 (Extended data Fig. 7) and were therefore treated as independent
388 variables. To analyze and visualize the relationship between the SDs calculated from variables
389 with different units, all values of SD were centered and scaled [z -scored SD = $(SD_{\text{observed}} -$
390 $SD_{\text{mean}})/SD_{\text{mean}}$]. The effects of CSV, SETUP and their interaction on among-laboratory SD was
391 tested with a mixed effects model using the “nlme” package²⁹ as suggested by Zuur et al.
392 (2009)³⁰. The statistical model with the lowest AIC for between-laboratory SD included the
393 response variable as a random factor as well as a “varIdent” weighting function to correct for
394 heteroscedasticity resulting for the variable-specific spread of the residuals (R syntax: “model=
395 lme (between-laboratory SD ~ CSV*SETUP, random=~1|variable, weights=varIdent (form =
396 ~1|variable))”). *A priori* planned contrasts between the CTR and the treatment levels with CSV
397 were performed using Welch’s t -tests on the z -scored normalized SDs.

398 Within-laboratory SDs were analyzed with two approaches. First, to allow for a direct
399 analytic and graphical comparison with the results for among-laboratory SD, we aggregated the
400 within-laboratory SDs by CSV and SETUP for each response variable ($n = 144$). A model
401 similar to the one we had used for among-laboratory SD was used to assess the impact of CSV
402 and SETUP. With a second approach we analyzed each response variable separately using mixed
403 effect models with “laboratory” as a random factor and a “varIdent” weighting function to
404 correct for heteroscedasticity resulting for the lab specific spread of the residuals (R syntax:
405 “model= lme (net legume effect ~ CSV*SETUP, random=~1|laboratory, weights=varIdent (form
406 = ~1|laboratory))”; $n = 84$; 14 laboratories \times 6 CSV treatments). As within-laboratory SD data
407 allowed us to account for the inherent collinearity of some of the response variables, we further
408 tested the impact of the CSV on the first and second principal components (PC1 and PC2)
409 derived from a principal component analysis (“prcomp” function in R) using scaled and centered

410 values from all twelve response variables. PC1 and PC2 were then analyzed with the same
411 mixed-effects model used to analyze the within-laboratory SDs from individual variables. The
412 relationship between within- and among-laboratory SD also was analyzed with a mixed-effects
413 model to test for effects on among-laboratory SD of CSV, within-laboratory SD, and their
414 interactions, and with the source response variable entering the model as a random factor.

415 As we observed a slight increase in reproducibility of within-laboratory SD for the CSV-
416 AB treatment, we further tested the effect of this treatment on the statistical power for detecting
417 the net legume effect in each individual laboratory. This analysis was performed with the
418 “power.anova.test” function as available in the “base” package. We computed the statistical
419 power of detecting a significant net legume effect (using a one-way ANOVA for the legume
420 treatment) for CTR and CSV-AB for each laboratory and response variable. This allowed us to
421 calculate the average statistical power for the two treatments and the extent to which the number
422 of replicate microcosms needs to be increased for CSV-AB to achieve the same statistical power
423 as for CTR.

424

425 **Additional references for methods**

- 426 27. Keuskamp, J. a., Dingemans, B. J. J., Lehtinen, T., Sarneel, J. M. & Hefting, M. M. Tea
427 Bag Index: a novel approach to collect uniform decomposition data across ecosystems.
428 *Methods Ecol. Evol.* 4, 1070–1075 (2013).
- 429 28. Team, R. C. R: A Language and Environment for Statistical Computing. (R Foundation
430 for Statistical Computing, 2016).
- 431 29. Pinheiro, J., Bates, D., DebRoy, S. & Sarkar, D. NLME: Linear and nonlinear mixed
432 effects models. R Packag. version 3.1-122, <http://CRAN.R-project.org/package=nlme> 1–

Milcu et al. 2016

433 336 (2016).

434 30. Zuur, A. F., Ieno, E. N., Walker, N., Saveliev, A. a & Smith, G. M. Mixed Effects Models

435 and Extension in Ecology with R. (2009). doi:10.1007/978-0-387-87458-

436 **Table 1 | Impact of experimental treatments on among- and within-laboratory SD.** Mixed
437 effects table output showing the impact of controlled systematic variation (CSV), experimental
438 SETUP (glasshouse vs. growth chamber) and their interaction on among- and within-laboratory
439 SD. Complementary analyses assessing the impact of experimental treatments on within-
440 laboratory SD can be found in Extended Data Table 3.

441

Source	Among-laboratory SD			Within-laboratory SD	
	df	F	<i>P</i>	F	<i>P</i>
CSV	5/121	4.58	0.001	2.27	0.052
SETUP	1/121	1195.70	<0.001	9.38	0.003
CSV×SETUP	5/121	2.48	0.035	0.27	0.926

442

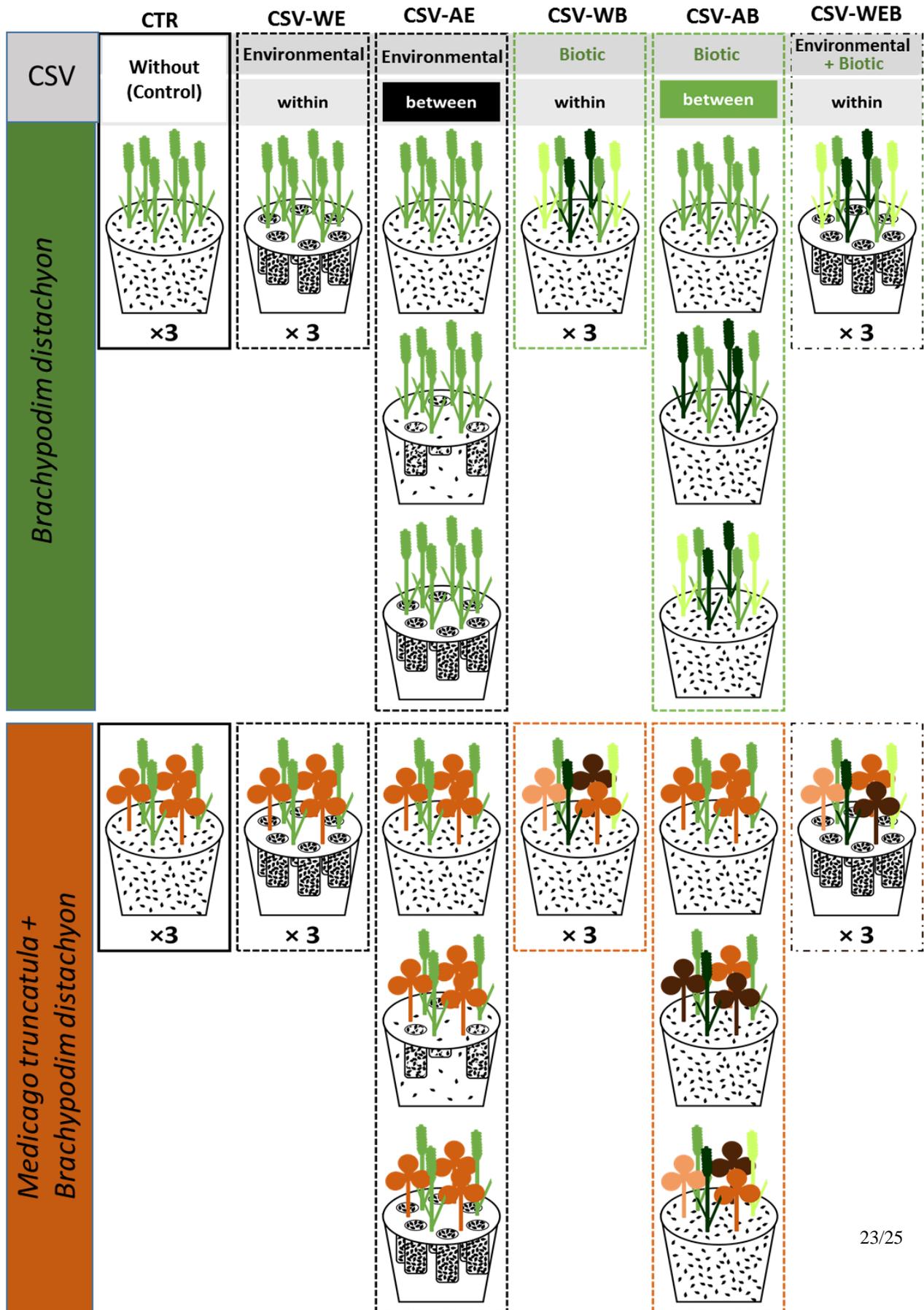
443

444 **Figures**

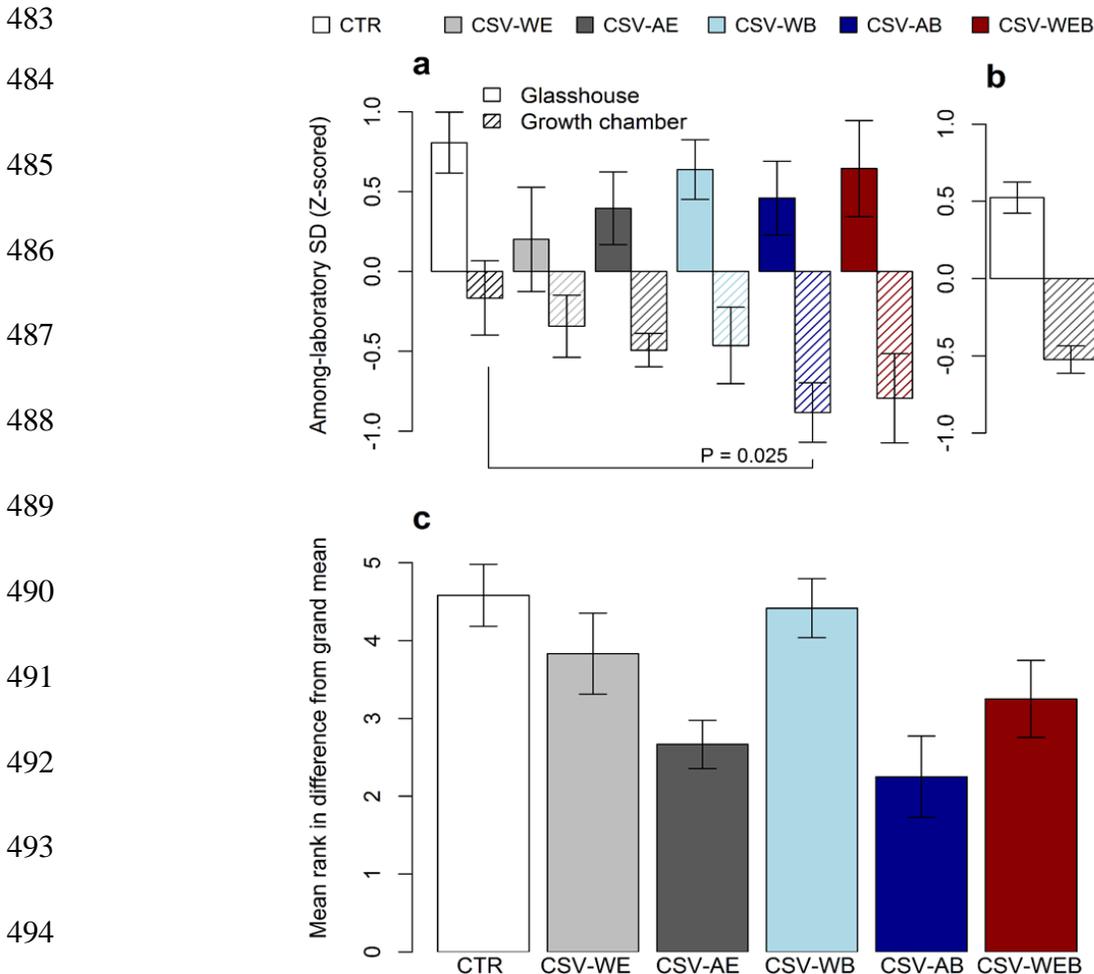
445 **Fig. 1 | Experimental design of one block.** Grass monocultures of *Brachypodium distachyon*
446 (green shades) and grass-legume mixtures with the legume *Medicago trunculata* (orange-brown
447 shades) were established in 14 laboratories, with the shades of green and orange-brown
448 representing distinct genotypes. Plants were established in a substrate with equal proportions of
449 sand (black spots) and soil (white), with the sand being either mixed with the soil or concentrated
450 in sand columns to induce environmental CSV. Combinations of three distinct genotypes were
451 used to establish biotic CSV. Alongside a control (CTR) with no controlled systematic variability
452 (CSV) and containing one genotype in a homogenized substrate (soil-sand mixture), five
453 different types of environmental or biotic CSV were used as treatments: 1) within-microcosm
454 environmental CSV (CSV-WE) achieved by spatially varying soil resource distribution through
455 the introduction of six sand patches into the soil; 2) among-microcosm environmental CSV
456 (CSV-AE), which varied the number of sand patches (none, three or six) among replicate
457 microcosms; 3) within-microcosm biological CSV (CSV-WB) that used three distinct genotypes
458 per species planted in homogenized soil in each microcosm; 4) among-microcosm biological
459 CSV (CSV-AB) that varied the number of genotypes (one, two or three) planted in homogenized
460 soil among replicate microcosms; and 6) both environmental and biotic CSV (CSV-WEB) within
461 microcosms that used six sand patches and three plant genotypes per species in each microcosm.
462 The “× 3” indicates that the same genotypic and sand composition was repeated in three
463 microcosms per block. The spatial arrangement of the microcosms in each block was re-
464 randomized every two weeks. The blocks represent two distinct chambers in growth chamber
465 setups, whereas in glasshouse setups the blocks represent two distinct growth benches in the
466 same glasshouse.

467

468



470 **Fig. 2 | Reproducibility as affected by experimental treatments. a**, Effects of CSV and
471 SETUP on among-laboratory SD. Lower values indicate enhanced reproducibility. Filled bars
472 and dashed bars represent glasshouse (n = 6) and growth chamber setups (n = 8), respectively.
473 The *P* value shown in **a** represents the result of an *a priori* planned contrast (Welch's *t*-test)
474 between CTR and CSV-AB treatment levels in the growth chamber setup **b**, Overall effect of
475 SETUP on among-laboratory SD. **c**, CSV effects on the mean rank difference from the grand
476 mean (resulted from pooling all 14 laboratories; see Extended Data Fig. 4). Lower difference
477 from the grand mean indicates increased reproducibility. A Kruskal-Wallis test on the ranked
478 differences from the grand mean for the twelve response variables found a significant effect of
479 CSV on the mean ranks ($\chi^2_{5,66} = 18.03$, $P = 0.003$), with the lowest mean rank difference (*i.e.*, the
480 closer to the grand mean) for the CSV-AB treatment and the highest mean rank difference for the
481 CTR treatment. Bars with error bars represent means ± 1 s.e.m., n = 12 (representing the z-scored
482 SD of the twelve measured response variables).



495 **Fig. 3 | Within-laboratory SD of the net legume effect.** **a**, Effects of CSV and SETUP on
496 within-laboratory SD. Filled bars and dashed bars represent glasshouse (n = 6) and growth
497 chamber setups (n = 8), respectively. **b**, Overall effect of SETUP on within-laboratory SD. **c**,
498 Relationship between within- and among-laboratory SD in growth chambers for CTR and CSV-
499 AB. **d**, Relationship between within- and among-laboratory SD in glasshouses for CTR and
500 CSV-AB. The twelve data points per CSV treatment (**c**, **d**) represent z-scored SD of the twelve
501 measured response variables. None of the regressions are significant at $P < 0.05$. Bars with
502 error bars represent means ± 1 s.e.m., n = 12. See Extended Data Fig. 6 for a graph presenting
503 the regression lines for all CSV treatments.
504

