Supplementary Information (SI) for Dee et al. "Clarifying the effect of biodiversity on productivity in natural ecosystems with longitudinal data and methods for causal inference."

- Data and Code Availability: The code for reproducing all analyses, figures, and tables in this
 study are available at https://github.com/LauraDee/NutNetCausalinf and released on Zenodo
 (DOI/10.5281/zenodo.7675340), which we refer to in this SI as our "project page." An
 RMarkdown tutorial on the methods can also be found on our Zenodo release and as a
- 9 Supplemental Data File with this publication.

10

Table of Contents

11	Supplemental Methods	3
12	S1. Glossary of Terms	3
13	S2. Directed Acyclic Causal Graph (Figure 1B)	7
14	S2a. Comparison to path diagrams and structural equation models	7
15	S3. Supplementary Methods: Data Description	
16	S3a. Measuring Biodiversity and Productivity	9
17	S3b. Data transformations	10
18	S4. Supplementary Methods: Main Design Estimator	
19	S4a. Review of Main Design from Methods section	13
20	S4b. Estimation procedure and implementation	15
21	S4c. Brief comparison of our design and aims to other study designs and aims	15
22	S5. Supplementary Methods: Extensions of the Main Design Estimator	16
23	S5a. Species evenness	17
24	S5b. Functional form	
25	S5c. Moderating effect of site-level species richness	
26	S5d. Moderating effect of site-level productivity	20
27 28	S6. Supplementary Methods: Robustness Checks to Assess Potential Threats to Inte Validity	ernal 23
29	S6a. Reverse causality: productivity causes species richness	

30	S6a.i. Blocked mechanism design	25
31	S6a.ii. Instrumental variable design	27
32	S6bi. Dynamic panel designs	31
33	S6b.ii. Design sensitivity to unobserved, plot-level confounding variables	37
34	S7. Comparison of Main Design to Common Designs in Ecology	41
35 36	S8. Supplementary Methods: Heterogeneous Effects of Rare, Non-rare, and Non-nati Species on Productivity	ve 45
37	S8a. Definitions and Measurement of groups in Figure 5	47
38	S8b. Statistical Analyses	49
39	S8c. Comparing the effect of species richness per group on productivity	50
40	S8ci. Sensitivity Analyses for species with unknown origins	51
41	S8cii. Sensitivity analyses using relative frequency as a metric for rarity	52
42	S8ciii. Sensitivity Analyses using different cut-offs for rare versus non-rare categories	54
43	S8d. Variation in each species group	55
44	Supplementary Discussion	58
45	S9. Nine Frequently Asked Questions (FAQ) about Dee et al	58
46	S10. Supplementary References	76
47		
48		
49		
50		
51		
52		
53		
54		
55		
56		
57		

58

59 Supplemental Methods

60 **S1. Glossary of Terms**

61 We provide brief definitions of terms, and we organize the terms logically rather than 62 alphabetically. For more details and background reading, see e.g., (1-7).

63 *Counterfactual (contrary to fact):* In the definition of a causal effect, a plot is assumed to have a 64 potential productivity outcome under each potential richness level; e.g., $P_i(R'')$ is the potential 65 productivity outcome when R = R'' and $P(R'_i)$ is the potential productivity outcome when R =66 $R' (R' \neq R'')$. But at any point in time, only one of those richness levels, and thus one of those 67 productivity values, will be observed. The other values are <u>counterfactual</u> values – i.e., the 68 productivity values that would have been observed had we instead observed the plot under the 69 other possible richness levels.

70 *Treatment*: Causal variables, like *R*, are often labeled "treatment variables" whether they are 71 manipulated by an experimenter or by nature. A change from one value to another is often 72 labeled a "treatment." See (6-8).

Average Treatment Effect (ATE): The ATE of biodiversity on productivity in plot *i* in year *t* is defined as $E[P_{it}(R'') - P_{it}(R')]$, where $E[P_{it}(R'')]$ is the expected productivity in plot *i* in year *t* when richness has the value R = R'' and $E[P_{it}(R')]$ is the expected productivity in plot *i* in year *t* when richness has the value $R = R' (R' \neq R'')$. The ATE is the average (or expected) causal effect of R on P for a randomly selected plot from the study population when biodiversity goes from R' to R''.

Directed Acyclic Causal Graph (DAG): A DAG is a visualization of qualitative causal
assumptions on which one relies for making causal claims from observable data (9). See Section
S2 for more information and the relationship between a DAG and a "path diagram."

82 Internal Validity: The extent to which a study design allows one to infer a causal relationship

83 from a correlation by ruling out rival explanations. For instance, are the changes in the

84 independent variable, X, causing a change in the dependent variable, Y, or can those changes in Y

85 be attributed to other causes?

External Validity: The extent to which inferences can be generalized (e.g., across sites, time
periods, contexts, or scales).

Construct validity: The extent to which an experimental treatment or statistical estimate matches
the phenomenon it intends to measure (10, 11) or the theory it intends to test.

90 Confounding Variables (or "confounders"): The term "confounding variable" describes

91 variables that are systematically correlated with the causal variable (e.g., biodiversity) and the

92 outcome variable (e.g., productivity), and thus can mask or mimic a causal effect. Confounding

93 variables are a potential source of *bias* in a study design.

94 *Bias and Hidden Bias:* An estimator is a rule or procedure for calculating an estimate of a causal 95 effect based on observed data. Bias is a property of an estimator: it captures the difference 96 between the estimator's expected value and the true value of the causal effect being estimated 97 (12). The phrase "hidden bias" (also called "unobserved heterogeneity") is often used to describe 98 the potential sources of bias in a study design (e.g., an omitted third variable that affects both 99 biodiversity and productivity). Hidden bias is thus a rival explanation for detecting or failing to 100 detect a correlation between a purported causal variable and its outcome using observable data 101 (reviewed in (13)). The goal of causal analysis is to choose data and a design so that an actual 102 causal effect would be visibly different from the most plausible hidden biases. Note: Sampling 103 variability ("noise" or "chance") is different from hidden bias. Sampling variability is a rival 104 explanation for observed relationships between variables, but it is not a source of bias. Sampling 105 variability declines with more data, whereas bias does not (14, 15). Sampling variability is 106 reflected in variable *I* in Figure 1B, whereas bias comes from variable *U*.

"Selection on Observables" Assumption: Informally, this assumption implies that confounding
variables that could introduce bias are known and observable to the researcher, so that statistical
bias can be eliminated (controlled, blocked) by conditioning strategies, such as regression,

110 matching or stratification methods. To read more, see (1).

111 *Fixed Effect:* Our use of the term "fixed effect" is drawn from the econometrics literature, where

112 it refers to the effect of a time-invariant attribute of the system (12); e.g., a plot-level fixed effect

113 is an attribute of the plot that is assumed to not change over the study period, such as topography

114 or historical patterns of land use. This use of the term "fixed effect" differs from how the term is 115 typically used in ecology, where the term often refers to the coefficient estimates of explanatory 116 variables in mixed (multi-level) modeling (e.g., (16)). Further confusing matters, the "random 117 effects" components of a mixed effects model, which describe categorical variables that are 118 assumed to be drawn from a normal distribution with zero mean, are often used to accomplish 119 the same goal as the "fixed effects" that we apply here (e.g. to remove spurious plot-level 120 effects). However, unlike "random effects," econometric fixed effects are not constrained to be 121 drawn from any predefined distribution. They are assumed to be fixed and estimable rather than 122 assumed to have a distribution (i.e., they are not part of the error term, as random effects are 123 assumed to be in multi-level modeling). Operationally, fixed effects in econometrics are simply 124 regression parameters describing categorical or dummy variables per study unit (e.g., in 125 experiments, a categorical fixed effect parameter per plot is often fit to control for differences 126 among plots that are not associated with the experimental treatment). Although this fixed-effect 127 estimation approach comes at the cost of reduced statistical power, it avoids the potential bias 128 that can arise when controlling for time-invariant variables using random effects (to read more, 129 see (3)). The use of random effects requires the assumption that the random effect is uncorrelated 130 with all of the covariates in the model (17). In an observational data set with any sort of 131 environmental gradient, that assumption is strong and not likely satisfied.

Mechanism: A mechanism is a variable that lies on the causal path between two other variables and mediates the causal effect of one of those variables on the other (*18*); shown as "M" in Figure 1 B (right panel). A mechanism can be viewed as an intermediate outcome of a causal variable; e.g., an increase in plot productivity causes a decrease in plot biodiversity by increasing the amount of shading in the plot – shading is the mechanism through which a change in productivity can cause a change in biodiversity.

Moderator: A moderator is a variable that lies <u>off</u> the causal path between two other variables but moderates the magnitude of a causal effect. A moderator is a source of heterogenous causal effects (18); for instance, the degree to which biodiversity affects productivity may depend on weather (e.g., precipitation or temperature) – weather can moderate the causal effect of biodiversity on productivity, but the change in biodiversity does not change weather conditions.

143 *Heterogeneous Treatment*: This concept goes by many names in the causal inference literature, including "multiple versions of the treatment", "treatment variation," "hidden versions", 144 145 "heterogeneous treatments," and "hidden treatments" ("hidden treatments" being used differently 146 from how ecologists have used the phrase in the past (19)). This issue can be viewed as a 147 challenge with construct validity: if you say that richness goes from 4 species to 8 species in plot 148 A, and I say that richness goes from 4 species to 8 species in plot B, we need to assess if we are 149 talking about the same change in the treatment variable. In this case, a richness change from 4 to 150 8 species can involve many combinations in species identity, even in experiments manipulating a 151 subset of all species in an ecosystem (e.g., are all four additional species native and not rare, or 152 are 2 of those additional species rare and 2 non-native?). The underlying idea is that, for a unit 153 of observation (e.g., a plot), there ought to be one potential outcome for each treatment value 154 (e.g., one potential productivity value for each richness value at a particular moment in time). If 155 there is not, we have multiple versions of the same treatment. In this way, species composition is 156 not a confounding variable, but a heterogenous treatment in the effect of richness on 157 productivity. Note that this concept of heterogeneous treatments is different from "heterogeneous 158 treatment effects," which simply means that not every unit responds the same way to a change in 159 the treatment variable (i.e., treatment effects are moderated by variables that differ across units in 160 the study population).

Instrumental variable: This term is defined in Figure 1B as variable Z, a variable that affects the
treatment variable (in our study, "species richness") but has no direct effect on the outcome
variable (in our study, "aboveground biomass") except through its effect on the treatment. In a
randomized experiment, the instrumental variable is the randomization procedure. To read more,
see (20–22) and for examples in ecology see (23, 24).

166 *Panel data*: Panel data are longitudinal data comprising repeated measures taken from a sample

167 of cases (e.g., plots, sites, regions). These data also called cross-sectional time series (25) or

168 longitudinal multilevel data. See section S3 for more detail.

169 Cross-sectional data: Data without repeated measures taken from a sample. Instead, one measure

170 is taken from each unit of analysis (e.g., plot). For example, when using cross-sectional data,

analyses of the effect of biodiversity on productivity have only one observation per plot.

172 *Errors versus residuals*: The error term in a regression model represents how the observations

173 differ from the true population. It is an unobservable and part of the true data generating process.

174 In contrast, residuals are an estimate of the unobservable error term, as the difference between

the regression line (predicted value) and the observed data points from the sample population.

176 Residuals cannot be used to assess potential bias in an estimation procedure, and thus this

177 distinction is important in our discussions of statistical bias.

178 S2. Directed Acyclic Causal Graph (Figure 1B)

179 Figure 1B (right panel) is known as a directed acyclic causal graph (DAG) and is a 180 visualization of qualitative causal assumptions (5, 9, 26–28). A DAG encodes knowledge and 181 beliefs about how a system works. The graphical relations depicted in the DAG encode causal 182 claims – not just representations of associations. A directed edge (e.g., $R \rightarrow P$) depicts a claim 183 about the results of many hypothetical experiments, whereby if every other variable represented 184 in the graph is held fixed, R and P will covary if R if manipulated, but not if P is manipulated 185 (note, time is implicit in the DAG, and a DAG assumes that one can isolate the effect of R on P 186 but does not imply that P can never affect R; another DAG may represent the reverse direction, P 187 \rightarrow R).

One key benefit of a DAG is that it makes transparent the assumptions on which one relies for making causal claims from observable data. A DAG therefore allows the researcher and the reader to better judge the credibility of the causal claims from a specific research design. Another way to view this benefit is that a causal graph helps identify the sources of variation in a causal variable and in its outcome, thereby emphasizing potential sources of bias that must be addressed in a research design and pointing to designs that can address these sources of bias (*1*).

194 S2a. Comparison to path diagrams and structural equation models

A DAG is like a "path diagram," which may be more familiar to ecologists and are often used in structural equation modeling (29, 30). Although not all path diagrams are DAGs, a DAG can be interpreted as a non-parametric structural equation model (SEM) (31) with no cycles (no double-headed arrows). In other words, an SEM can be a DAG, but an SEM could also contain both cyclic and directed cycles (not a DAG).

200 In practice, however, SEMs, when used for causal claims, rely on conditioning on observable 201 confounding characteristics to eliminate non-causal dependencies between two variables ('the selection on observables assumption' - see S1 Glossary). However, DAGs emphasize also making transparent assumptions about unobservable confounding variables. In contrast to the common practice of SEMs in ecology, our design can eliminate unobserved confounders (section S5). Nevertheless, SEMs, as typically implemented in ecology, have advantages over our design in cases where a researcher believes that all important confounders can be observed and controlled within the SEM: in those cases, SEMs can be more efficient (i.e., higher statistical power) and they expand the scope of analyses that can be performed with a single data set and estimation strategy.

S3. Supplementary Methods: Data Description

We analyze panel data from grasslands around the world in the Nutrient Network (32, 33), which includes mesic grasslands and prairies, savanna, desert grasslands, montane meadows, old fields, and alpine tundra. We use data from 43 sites with unmanipulated plots with at least 5 years of data in the period 2007-2017 (see Table S1). Unmanipulated plots are control plots in the nutrient addition experiments of the Network, meaning they receive no additional nutrients. Unprocessed data versions were 'full-cover-09-April-2018.csv', and 'comb-by-plot-clim-soil-diversity-09-Apr-2018.csv' from the Nutrient Network. All R scripts to process data and create derived data is available at the project page (DOI/10.5281/zenodo.7675340).

Table S1. Information on unmanipulated control plots from the Nutrient Network. The table shows the number of plots with data by year and site. All plots in the analysis have at least 5 years of data between 2007-2017. The dataset includes sites from 11 countries and 5 continents (North America, Australia, Europe, South America, and Africa).

2012 2013 Site Name Bogong _ _ **Boulder South Campus** ----_ Bunchgrass (Andrews LTER) Burrawan _ _ _ Cedar Creek LTER _ **Cedar Point Biological Station CEREEP** - Ecotron IDF _ _ _ _ _ **Chichaqua Bottoms** _ _ -Companhia das Lezirias -_ _

	-	-	-	r	1	-	r	r	1	1	1
Cowichan	3	3	3	3	3	3	3	3	3	-	-
Doane College Spring Creek Prairie	-	-	-	-	-	2	2	2	2	2	2
Duke Forest	3	3	3	3	3	-	-	-	-	-	-
Elliott Chaparral	-	-	3	3	3	3	3	2	3	3	-
Fruebuel	-	3	3	3	3	3	3	-	3	-	-
Hall's Prairie	3	3	3	3	3	3	3	3	-	-	-
Hart Mountain	3	3	3	3	3	3	-	-	-	-	-
Heronsbrook (Silwood Park)	-	3	3	3	3	3	-	-	-	-	-
Hopland REC	3	3	3	3	3	3	3	3	3	3	3
Kinypanial	-	-	3	3	3	3	3	3	3	-	-
Koffler Scientific Reserve, Joker's Hill	-	-	-	9	9	9	9	9	9	9	9
Konza LTER	3	3	3	3	3	3	-	3	3	-	-
Lancaster	-	3	3	3	3	-	-	3	3	3	-
Lookout (Andrews LTER)	3	2	3	3	3	3	3	3	3	3	3
Mar Chiquita	-	-	-	-	3	3	3	3	3	3	-
Mclaughlin UCNRS	3	3	3	3	3	3	3	3	3	3	3
Mt. Caroline	-	4	4	4	4	4	4	4	4	4	-
Papenburg	1	1	1	1	1	1	1	-	-	-	-
Rookery (Silwood Park)	-	3	3	3	3	3	-	-	-	-	-
Sagehen Creek UCNRS	3	3	3	3	3	3	3	-	-	-	-
Saline Experimental Range	-	3	3	3	3	3	3	3	3	-	-
Savannah River	2	2	2	2	2	2	-	-	-	-	-
Sedgwick Reserve UCNRS	6	6	6	6	6	6	6	6	6	6	6
Serengeti	-	3	3	3	3	3	-	-	-	-	-
Sevilleta LTER	5	5	5	5	5	5	5	5	-	-	-
Sheep Experimental Station	4	4	4	4	4	4	-	-	-	4	-
Shortgrass Steppe LTER	3	3	3	3	3	3	3	3	3	3	-
Sierra Foothills REC	3	5	5	5	5	5	5	5	5	5	5
Smith Prairie	3	3	3	3	3	3	-	-	3	-	-
Spindletop	3	3	3	3	3	3	3	3	3	3	3
Temple	4	4	4	4	4	4	4	4	4	4	-
Trelease	-	-	3	3	3	3	3	-	-	-	-
Ukulinga	-	-	6	6	6	6	6	6	6	6	-
Val Mustair	-	3	3	3	3	3	3	3	3	3	-

226

227

228 S3a. Measuring Biodiversity and Productivity

229 To measure productivity, we use plant above-ground live mass (biomass) as in **Figure S1**.

230 Biomass production supports many ecosystem processes and services and this measure of

231 productivity has been widely used in addressing the relationship between diversity and

productivity with observational data (e.g., (34–36)) and in many grassland experiments

233 (reviewed in (37–39)). For herbaceous vegetation, above-ground live biomass provides a

reasonable estimate of primary productivity (40).

235 Live above ground biomass is measured in the Nutrient Network dataset with the following 236 procedure. In each 5m x 5m unmanipulated control plot, a permanently marked, randomly 237 located, 1m x 1m subplot is sampled annually at peak biomass for species composition. Visual 238 cover estimates are made to the nearest 1% for every species contained within (or over-hanging) 239 the subplot and used to calculate species diversity metrics (richness, evenness). Biomass samples are collected from two 1 m x 0.1 m strips (totaling 0.2 m²) located adjacent to the 1m² cover 240 241 subplot. All vegetation from plants rooted within these strips is clipped at ground level. The 242 location of the biomass plots changes yearly to avoid repeat sampling previously clipped areas. 243 Biomass is dried at 60°C to constant mass and weighed to the nearest 0.01g. Multiplying weights 244 by five generates a gram per square meter value for productivity.

245 To make our study comparable with previous studies, we measure biodiversity as species 246 richness, the number of species in a plot in each year (Figures S1 and S2). We also consider 247 (Table S2) analyses that include species evenness, measured as the degree of similarity in 248 abundance between species within a community (41), and analyses that measure diversity with 249 Simpson's Diversity. We calculate Simpson diversity as the inverse Simpson index: 1/D where D 250 = sum (p_i^2) and p_i is the proportional index of each species *i* in a plot. The evenness variable is 251 $H/\log(S)$ where H is the Shannon Index and S is the number of species in a plot. $H = - \sup(p_i)$ 252 * $\ln(p_i)$). We calculate each metric using all the vegan package in R (42).

In this dataset, we started with 1291 potential plot-year observations but had to drop 2 plotyear observations because of missing richness values and another 58 observations because of missing productivity values.

256 S3b. Data transformations

257 Prior to estimating the effect of diversity on productivity, we transform our productivity
258 variable (live biomass) and our diversity variables (richness, evenness, Simpson's index) by
259 taking the natural logarithm of each plot-level measure. This transformation has several

260 advantages, which are all related in a statistical sense. First, both productivity and richness are 261 strictly positive variables that exhibit right-skewed distributions (see Figure S1). Transforming 262 by the natural logarithm reduces the skew of these variables, improving statistical efficiency (i.e., 263 improves the precision of our estimates). Second, in an ecological sense, it is reasonable to 264 assume that going from 2 to 4 species will on average have a bigger effect on productivity than 265 going from 18 to 20 species but may have a similar proportional effect on average as a change 266 from 10 to 20 species would have. In other words, the natural logarithm transformation makes 267 sense in situations when it is better to compare relative changes rather than absolute changes. In 268 other words, instead of assuming that P increases as a constant function of R, we assume that P

269 270

Figure S1. Plot-level (a & b) and site-level (c & d) species richness and productivity (aboveground live mass) between 2007-2017. (a) shows the levels of richness and productivity in all plots and (b) shows the log of richness and productivity in all plots. For comparison, (c) shows the average levels of richness and productivity across the 43 sites (see Table S1) and (d) the log of richness and productivity across the same sites.



- 277 increases as a relative function to the current level of P as a function of R. Another way to say 278 the same thing is that in a graph with richness on the horizontal axis and productivity on vertical 279 axis, a straight line will not the best description of the relationship. Third, the coefficient on 280 richness in our log-log specification has a well-define interpretation, which is a valuable trait; for 281 most readers, a single coefficient is more accessible and easier to evaluate than a non-linear 282 surface. In this SI (section S6), we also present the estimated effects of richness on productivity 283 in levels (i.e., no transformation), including quadratic and cubic specifications that permit the 284 estimated relationship to be non-linear.
- 285





289

292 **S4. Supplementary Methods: Main Design Estimator**

293 S4a. Review of Main Design from Methods section

We present the details on the Main Design in the main text *Methods* section. Here, we elaborate on the estimation procedure used to implement the regression model for the Main Design. In our study design, an observation comes from a plot *p* located within a site *s* in a year *t*. Recall that, to eliminate the confounding effects of time-invariant plot attributes (δ_{ps}) and time-varying site attributes (μ_{st}), we estimate an equation of the following form:

299

300
$$\ln LiveMass_{pst} = \beta \ln Richness_{pst} + \delta_{ps} + \mu_{st} + \varepsilon_{pst}$$
 (S1)

301 Given that we have a ln-ln specification, β can be interpreted as an elasticity: the expected 302 percent change in productivity given a one percent change in richness. In the economics 303 literature, the <u>time-invariant</u> plot attributes (δ_{ps}) would be called "plot-level fixed effects." Note 304 that fixed effects have a different meaning in economics than in ecology (see S1 Glossary). In economics, including δ_{ps} is said to control for "unobserved heterogeneity" across plots that can 305 306 be a potential source of bias. Note that δ_{ps} is not part of the error term, as it would be in mixed 307 (multi-level) models (see Section S7). Rather, it is a parameter to be estimated, just like β (i.e., β and δ_{ps} are assumed to be fixed and estimable, rather than assumed to follow a distribution). 308 309 Time-invariant site attributes are not explicitly included in the equation because they are 310 subsumed into the time-invariant plot attributes (i.e., plots are nested within sites and so fixed 311 site attributes are controlled via fixed plot attributes).

The <u>time-varying</u> site attributes (μ_{st}) are modeled in a fully flexible way that allows a yearspecific effect for each site (in the estimation, an indicator for each year is interacted with an indicator for each site). Explicitly estimating μ_{st} flexibly controls for confounding variation due to conditions at a site that vary from year to year, namely weather (e.g., temperature, precipitation), drought events, grazing, surrounding management, or other site-level attributes

that change through time. In other words, this variable captures all year-specific conditions

318 experienced by every plot at a given site.

319 The term ε_{pst} is a time-varying random error term at the plot level, assumed to have mean 320 zero and no correlation with ln *Richness*, i.e., it corresponds to I_{pst} in Figure 1B. Errors at a 321 given plot (ε_{pst}) may be serially correlated (i.e., temporally dependent even after conditioning on 322 richness and site-by-year effects), and thus we cluster the standard errors at the plot level (43).

323 Our clustered estimation of the variance allows for arbitrary serial correlation within each plot,

given site may also be correlated (even after conditioning on site-by-year effects) and thus, as a

- 324 as well as heteroskedasticity across plots (17, 44)). See our project page for code. Errors at a
- 325 given site may a
- 326 robustness check, we also estimate standard errors clustered at the site level (Table S2).

327 To sum up our design, we are asserting that, after controlling for time-invariant plot attributes 328 that are correlated with richness and productivity, and time-varying site attributes that are 329 correlated with richness and productivity, the remaining variation in richness in a plot is "as if 330 randomly assigned," independently across time. In other words, the remaining variation in 331 richness is driven by variables that have no link to productivity other than through their effect on 332 richness (i.e., Z_{pst} in Fig. 1). If our assumption is correct, we can give a causal interpretation to 333 the estimate of β . In section S7, we describe how we explore the sensitivity of our causal 334 interpretation to violations of this assumption. The estimated effect of richness on productivity is 335 reported in Figure 2 and Table S2.

336 As we noted in the Methods section of the main text, we seek to estimate the average causal 337 response of an incremental change in R across all plots (i.e., the average effect across all possible 338 one-unit changes). Recently, scholars have identified a potential form of misspecification bias 339 that may arise when using models like Equation (S1) to estimate this average causal response 340 when treatments are multi-valued and time-varying and the average causal response is 341 heterogeneous across time or treatment values (45). Specifically, the regression estimator applies 342 weights to all of the richness contrasts and year contrasts in the data and these weights can, in 343 some rare cases, be negative. In the presence of heterogeneous average causal responses, such 344 weights could overweight or underweight specific contrasts in a way that would create bias. This 345 bias can arise when (1) the average causal response wanes or matures over the panel when 346 treatment values change and remain at their new value across years (i.e., when they move to an 347 absorbing state); and (2) the distribution of treatment values is highly non-normal and the 348 average causal responses at extreme values of richness differ from the average causal responses 349 for values in the middle of the distribution of richness values. There is no theoretical basis for the 350 first source and the richness values in our panel data are well approximated by a normal 351 distribution. Thus, we do not believe that this form of specification bias is a potential problem in 352 our Main Design.

353 *S4b. Estimation procedure and implementation*

354 To ensure transparent, reproducible results among a wide range of scientists, we estimated 355 Equation (S1) in two software programs and by multiple coauthors for reproducibility. We used 356 the "reghdfe" command in STATA (v.16) (the "xtreg, fe" command yields the same estimates) 357 and the "feols" command in R in the fixest package (v. 0.8.2). While there are other packages in 358 R to execute this estimator, e.g., using the "felm" command in "lfe" (v 2.8-5) (46), we opt to use 359 "fixest" because the standard error estimation matches STATA and yields a more conservative 360 estimate based on the finite sample degrees of freedom correction for multi-way clusters. There 361 is no consensus on the "correct" finite sample degrees of freedom correction for these models, so 362 we opted for the more conservative option that results in larger standard errors.

363 *S4c. Brief comparison of our design and aims to other study designs and aims*

364 In our Main Design, our notion of causality and our approach differ from predictive (best-365 fitting) modeling approaches that use time-series data, such as convergent cross mapping designs 366 (47). Our "intervention-based" notion of causality (9, 48, 49) is what experimentalists have in 367 mind when they make causal claims (50). Furthermore, our model is not intended to be the best 368 *predictive* model of productivity;¹ i.e., the best model for predicting the level of productivity in plots outside of our sample (51). In fact, the best predictive model of productivity may not even 369 370 include richness as a variable. However, we are not interested in out-of-sample prediction of 371 plot-level productivity. In contrast, our goal is to infer the causal effect of richness on 372 productivity. But if one wanted to do such prediction, our design would pose challenges because 373 our estimates are conditional on the sample; the plot-level fixed effects are not assumed to have a 374 distribution (like they would in a mixed model), but rather are instead treated as fixed and 375 estimable.

The approach in our Main Design also differs from a mixed-effect modeling approach more common in Ecology (*16*). While ecologists who are familiar with multi-level modeling may wonder why, given our data are comprised of plots nested within sites and annual observations nested within plots, we do not use this multi-level modeling approach as our Main Design. We opt to use our design, rather than a mixed effect model, because our approach makes weaker and

¹ This post by Paul Allison (2014) explains key differences in evaluating multivariate regression models for the aim of prediction versus causal inference: https://statisticalhorizons.com/prediction-vs-causation-in-regression-analysis

381 more plausible assumptions for our data context and question, compared to a mixed effect model.

- 382 A full explanation is beyond the scope of this paper, but the main reason has two parts, which is
- 383 laid out in more detail in Section S7. First, without more variable transformations, the multi-level
- 384 modeling approach does not easily lend itself to controlling for as many unobservable sources of
- 385 confounding as can be done in our estimator. Second, clustering our standard errors at the plot
- 386 provides the same benefits that multi-level modeling does when estimating the variance-
- 387 covariance matrix in the presence of intra-site correlations among plots (44). For more detail and
- 388 discussion, see Section S7. Comparison of Main Design to Common Design in Ecology (a.k.a.,
- 389 multi-level modeling, hierarchical modeling, random effects modeling, mixed effects modeling,
- 390 or variance components modeling).
- 391

Table S2. Supplementary results on different variations of the main design. Column (1)
presents the results presented in the main text for the Main Design. The negative effect of ln of
species richness on ln of productivity holds when clustering standard errors at the site level
(column 2), when controlling for species evenness (column 3 & 6), and when using other
measures of biodiversity (Simpson's Diversity – column 4) as well as the lagged effect of species
richness in the prior year (ln *SpeciesRichnesst-1*) (columns 5 & 6). The estimated effect in column
(1) is plotted in Figure 2.

	(1)	(2)	(3)	(4)	(5)	(6)
ln(SR)	-0.2418 ***	-0.2418 **	-0.2237***		-0.2185**	-0.2057**
	(0.0854)	(0.0892)	(0.0851)		(0.0939)	(0.0948)
[-0.4	40902; -0.0743]	[-0.4165; -0.0670]	[-0.3905; -0.0568]		[-0.4024; -0.0345]	[-0.3914; -0.0199]
ihs(Evenness)			-0.1864			-0.1450
			(0.2122)			(0.2387)
			[-0.6022; 0.2294]			[-0.6128; 0.3228]
ln(Simpson)			-	0.1701 **		
			(0	.0679)		
			[-0.30	31; -0.0370]		
ln(lagged SR _{t-}	1)		_	_	-0.0146	-0.0096
					(0.0905)	(0.0903)
					[-0.1919; 0.1627]	[-0.1866; 0.1675]
Num. obs.	1231	1231	1231	1231	1093	1093
Num. plots	151	151	151	151	151	151
\mathbb{R}^{2} (full mod	el) 0.87	0.87	0.87	0.87	0.87	0.87

421 Robust Standard errors in parentheses, clustered at plot level in column 1 and clustered at site level in column 2.

422 423

424 S5. Supplementary Methods: Extensions of the Main Design Estimator

In this section, we present supplementary results, including variations in the specification ofour Main Design shown in Equation S1 and Equation 2 in the main text.

427 S5a. Species evenness

428 The negative estimated effect of richness on productivity (Figure 2) could reflect changes in 429 evenness, which may covary with changes in richness. In our sample, species richness varies a 430 lot over time within plots, but evenness does not (Figure S3). Thus, we do not suspect that failing 431 to include species evenness in Equation (S1) is a source of bias, but we nevertheless re-estimate 432 the equation after adding a measure of evenness. We transform the evenness variable with an 433 inverse hyperbolic sine (IHS) transformation, which has an effect similar to the natural logarithm 434 transformation but, unlike the natural logarithm transformation, is appropriate for variables, like 435 changes in evenness, that have many zero values (the natural logarithm of zero is undefined; 436 (52)).

437 After accounting for evenness, the negative relationship between richness and productivity 438 remains unchanged (Table S2), implying that the estimated effect of species richness in the first 439 column comes from changes in species richness rather than evenness. The estimated coefficient 440 on evenness is imprecisely estimated (Table S2) (i.e., large standard errors). Thus, anyone 441 interested in estimating its causal effect on productivity would not be able to draw precise 442 inferences from our data. This imprecision highlights that the strength of our design – its ability 443 to leverage change in diversity over time within plots to isolate the causal effect of diversity – 444 can be a weakness when the variable of interest does not change much over time within plots. If 445 this lack of temporal variation is common in many non-experimental contexts, experimental 446 designs varying evenness (e.g. (53)) may be the only way to obtain precise estimates of the role 447 of evenness on productivity or other ecosystem functions.

448



449

450 Figure S3. Year-to-year variation in species evenness per plot for the dataset described in Table451 S1.

452

453 *S5b. Functional form*

454 The effect of species richness on productivity could vary by the magnitude of the change in 455 the number of species. To detect this non-linearity, we estimate a quadratic specification of our 456 equation in which the variables are not log transformed (Table S3: 'Quadratic' column). We 457 detect evidence that the negative effect of richness on productivity becomes smaller as richness 458 increases, i.e., the coefficient on the non-squared term is negative and the coefficient on the 459 squared term is positive. Species richness ranges from 1 to 37 species in our dataset, which 460 determines the range over which there is a positive or negative effect. With the quadratic 461 specification, we find that the estimated effect or richness on productivity only turns positive in 462 plots over 31 species, which represents only 14 observations and 1.14% of the data (see Table S3 - 'Quadratic' column). For completeness, we also present estimates with the linear specification 463 464 and untransformed variables (level-level) and estimates with only productivity log-transformed 465 (log-level) (Table S3). We also estimated a cubic specification, and the estimated cubic term was 466 quantitatively and statistically indistinguishable from zero (see code on our project page, 467 DOI/10.5281/zenodo.7675340).

- 468
- 469

Table S3. Estimates of the effect of species richness on productivity *P* under changes in model specification for the functional form of this relationship. The columns compare the estimates from main log-log model (in column 1) to estimates from models with (2) the *ln* of productivity *P* but untransformed richness (i.e., richness in levels), (3) untransformed richness and productivity *P*, and (4) untransformed richness and productivity *P* with a quadratic term for richness. Standard errors, clustered at plot level, are in parentheses. 95% CIs are in brackets. All models include plot and site-by-year fixed effects as in Equation (S1).

	$(1)\ln(P)$	(2) $\ln(P)$	(3) <i>P</i>	(4) P (Quadratic)
ln(Richness)	-0.2418**	**		
	(0.0854)			
[-	0.4092, -0.0	0743]		
Richness		-0.0147*	-1.916	-16.83**
		(0.0080)	(2.919)	(6.623)
		[-0.0304, 0.0009]	[-7.637, 3.805]	[-29.81, -3.850]
Richness ²				0.5602**
				(0.2162)
				[0.1364, 0.9840]
Num. obs.	1231	1231	1231	1231
Num. plots	151	151	151	151
R^2 (full model)	0.867	0.865	0.83	0.83

494 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

495

496 S5c. Moderating effect of site-level species richness

We tested the hypothesis that site-level richness moderates the effect of plot level richness on productivity. This hypothesis is motivated by the observation that one potential reason for the observed negative effect of plot richness on plot-level productivity is that the "best performing" species enter a plot and win by becoming more productive (*54*) and thus causes a decline in richness to lead to an increase in productivity. This relies on having a large pool of species at a site that can colonize to take over during specific years.

We found no evidence supporting the hypothesis that the effect of plot-level species richness on plot-level live biomass depended on the levels of the observed site richness across all years or per year, nor on the numbers of introduced and native species at the site level (Table S4). We recommend that future research could test this hypotheses using data that includes direct observation of dispersal patterns (e.g., (54)), which are not available for this dataset at present. 508 Table S4. No dependence of the plot-level species richness (SR) effect on site-level species

509 richness (site SR) characteristics on productivity. All estimates are on a ln-ln scale. We

510 consider several site-level species richness measures, including: site-level richness across all

511 years (Site SR): count of all unique taxa ever observed across all plots in all years at that site), 512 the site-level richness per year (Site SR per year): count of unique taxa observed across all plots

512 the site-level richness per year (Site SR per year): count of unique taxa observed across all plots 513 at the site in that year), the count of all unique introduced taxa at the site (Site Introduced SR)

and the count of all unique native taxa at the site (Site Native SR). Interactions are indicated by a

515 "x." If anything, we find evidence that controlling for site level richness variables makes the

516 estimated effect of the log of plot species richness on log productivity more negative. To see

517 95% confidence intervals as well, see the project page (TableS4_R_CI.tex).

			Mode	<u>l:</u>	
	Total Si	te SR	Site Introduced SR	Site SR by year	Site Native SR
ln (SR)	-0.4	398 **	-0.2753**	-0.3834 *	-0.3184**
	(0.2	2155)	(0.1197)	(0.1998)	(0.1494)
ln(SR) x Site SR	0.0	0026			
	(0.	0027)			
ln(SR) x Site Int	roduced SR		0.0018		
			(0.0061)		
ln(SR) x Site SR	l per year			0.0041	
				(0.0050)	
ln(SR) x Site Na	tive SR				0.0021
					(0.0031)
Num. obs.	1231	123	1 1231		1231
Num. plots	151	15	1 151		151
R^2 (full model)	0.87	0.8	0.87		0.87

537 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

538 Robust Standard errors in parentheses (clustered at plot level).

539 S5d. Moderating effect of site-level productivity

A recent study by Wang et al. (55) found that the effect of biodiversity on productivity was moderated by the average level of productivity at a site, meaning that effect of biodiversity on productivity differed between high versus low productivity sites. In response, we test whether our results are altered by considering site-level productivity as a moderator. We consider sitelevel productivity in four ways: two using continuous variables and two using the cut-offs for high, medium, and low productivity classifications from Wang et al (55). Wang et al (55), however, used cross-sectional analyses to estimate this effect. Here, we can measure site-level

- 547 productivity in two ways: average over the entire time series ('Average prod. per site'), and site-
- 548 level productivity per year ('Average Prod per site & year'). In all calculations of site-level
- 549 productivity, we include the average productivity for the unmanipulated (control) plots, but not
- 550 the experimental plots at the Nutrient Network experimental sites.
- 551 We expand and estimate our main model in (equation S1) adding an interaction term between
- 552 In *Richness_{pst}* and *Avg Productivity* (Table S5). Next, Wang et al (55) categorize sites as high,
- medium, or low productivity in the 151 grids in HerbDivNet data based on mean productivity in
- 554 a grid: low between 30.18-238.73 (g/m²), medium between 239.67-409.69 (g/m²), and high
- between 414.29-1382.42 (g/m^2). Our productivity, in terms of live aboveground biomass at the
- site average across years, ranged from 62.48 to 1124.27 g/m^2; whereas the average site-level
- productivity per year ranged from 5.372 to 1609 g/m^2. Thus, to be comparable to Wang et al
- 558 (55), our groups were classified as: low below 239.67, high over 414.29, and the rest of sites as
- 559 medium productivity. We adopt these cut-offs and rerun the models interacting the ln *Richness*_{pst}
- 560 with the *Productivity_Group* (see Table S5 for details).

561 **Table S5. Estimating the moderating effect of site-level productivity on the effect of plot**-

level species richness (SR) on productivity, using continuous measures of productivity. As
 moderators, we consider the average site-level productivity per year (column 1) and the average
 site-level productivity across years (column 2). Interactions with plot-level species richness per
 year are indicated by a "x." To see 95% confidence intervals as well, see the project page
 (TableS5_R_CI.tex). All estimates are on a ln-ln scale.

		Moo	lel:
		(1)	(2)
n(SR)		-0.3305 **	-0.3532**
		(0.1609)	(0.1635)
n(SR) x Ave. Site	Prod. Per Yr.	0.0003	
		(0.0004)	
n(SR) x Ave. Site	Prod.		0.0004
			(0.0004)
Num. obs.	1231		1231
Num. plots	151		151
R^2 (full model)	0.87		0.87

582 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

583 Robust Standard errors in parentheses (clustered at plot level).

Across these analyses, we can detect no moderating effect of site-level productivity, based on estimated coefficients and their SEs in Table S5 and S6. See R code to reproduce analyses at the project page.

587 Table S6. Estimating the moderating effect of site-level productivity on the effect of plot-588 level species richness (SR) on productivity, using categorical groups of high, medium, low 589 productivity based on Wang et al. We interact plot-level richness with each productivity group; 590 interactions are indicated with an 'x' in the results table. Productivity groups were determined as 591 follows. Wang et al (55) use a single year of data; to mirror this measure, we use an average 592 productivity per site across years (column 2). In contrast to Wang et al (55), we also interact 593 plot-level richness with an average productivity per site per year (column 1). To see 95% 594 confidence intervals as well, see the project page (TableS6 R CLtex). All estimates are on a ln-595 ln scale. 596 597 (1)(2)598 Average Prod. per site & year Average Prod. per site 599 _____ 600 ln(SR) -0.2090* -0.1946 601 (0.1070)(0.1395)602 0.1075 ln(SR)x ProdGroupMedium 603 (0.1155)604 ln(SR) x ProdGroupHigh -0.0718 605 (0.1740)606 ln(SR) x ProdGroup:WangCutoffsMedium -0.1358 607 (0.1985)608 ln(SR) x ProdGroup:WangCutoffsHigh 0.0623 609 (0.2123)610 _____ 611 Num. obs. 1214 1231 612 R^2 (full model) 0.77 0.77 613 Num. plots 151 151 614 615 Signif. Codes: ***: 0.01, **: 0.05, *:0.1 616 Robust Standard errors in parentheses (clustered at plot level).

617 618

010

619

620 S6. Supplementary Methods: Robustness Checks to Assess Potential Threats 621 to Internal Validity

622 In our Main Design (Section S4), the key, untestable assumption for drawing a causal 623 inference from our estimator is that, after controlling for time-invariant plot confounders and 624 time-varying site confounders, the remaining variation in richness in a plot is "as if randomly 625 assigned," independently across time. In the main text (Figure 3), we consider potential 626 violations of this assumption and the implications for our inferences. More specifically, we 627 conduct a series of analyses that rely on alternative assumptions for causal inference. As noted in 628 the main text, the results are consistent across all approaches. Here, we describe these 629 approaches in more detail.

630 First, we explore potential violations in our assumption that the effect we are estimating goes 631 from richness to productivity, and not the other way around (in Section 6a). Because species 632 richness and biomass measures are taken simultaneously each year, as they typically are in many 633 ecological data sets, we cannot rely on temporal sequencing of the data to rule out reverse 634 causality. To address this potential threat to causal inference in our design, and in the process 635 also address potential bias from unobserved, time-varying plot attributes, we take two 636 approaches: (a) we posit a mechanism through which productivity affects richness -i.e., shading 637 (based on (56)) – and then block this mechanism and evaluate the change in our estimated effect 638 of richness on productivity (Section 6a.i); and (b) as an alternative to our main estimator 639 (Equation S1), we use an estimator that can estimate the effect of richness on biomass for a 640 subsample of the observations for which we can more credibly argue that the direction of 641 causality goes from richness to productivity (Section 6a.ii).

642 After assessing the potential threat to inference from reverse causality, we then explore 643 potential violations in our assumption that there are no time-varying plot attributes that are 644 systematically correlated with richness and productivity. To do this, we take two approaches. 645 First, we explore violations in our assumption that prior productivity does not influence current 646 richness, an effect that could be mediated by prior species richness (i.e., reverse causality in prior 647 year) or by other dynamic mediators. To address this potential source of bias from a plot-level, 648 time-varying confounder, we use two alternative estimators (Section S7b.i) that replace the plot 649 "fixed effects" with lagged productivity (i.e., lagged dependent variable estimators). Second, we 650 take a more general approach to quantifying our uncertainty about the potential bias from a timevarying, plot-level confounders. We create bounds our estimated targeted causal effect by
assuming that there are time-varying plot attributes that are systematically correlated with
richness and productivity. Specifically, we explore how our estimated effect would change if
there were an unobserved confounder that was negatively correlated with richness and positively
correlated with productivity (i.e., a source of bias that yields a spurious negative causal
relationship between richness and productivity in our design; *Section S7b.ii*).

657 S6a. Reverse causality: productivity causes species richness

658 For changes in richness to cause changes in productivity, changes in richness must occur 659 prior to changes in productivity. However, as in most experimental and observational studies on 660 the relationship between diversity and productivity, the Nutrient Network data on diversity and 661 productivity are collected at the same time each year. In the absence of high-resolution temporal 662 data (e.g., daily), we must make additional assumptions and run additional tests to rule out 663 reverse causality. If productivity were to negatively affect richness, that causal relationship could 664 mask a positive relationship of R on P in our design. In other words, our estimated coefficient of 665 β in Equation (S1) may reflect a causal relationship that runs from productivity to richness, rather than the other way around (i.e., it reflects a causal graph with a directed edge that flows 666 667 from *P* to *R* instead from *R* to *P*).

To illustrate the problem caused by reverse causality, we create a new causal graph in Figure 668 669 S4 (this graph is not acyclic because it has bi-directional arrows between two variables). We use 670 the notation from Equation (S1) but suppress the plot subscript p and add a subscript t-1 for time 671 lagged one period. In this new graph, we assume that bias from regressing P_t on R_t is not coming 672 from unobserved confounders, but rather from simultaneous causal relationships in which P_t and 673 R_t directly cause one another. If we regress P_t on R_t and the estimated coefficient β is less than 0, 674 a critic of our analysis could argue that the true β is greater than zero but masked because $\alpha < 0$ 675 and $|\alpha| > |\beta|$. If P_t affects R_t, ε_t is necessarily correlated with R_t, which is a violation of our 676 assumption that, after controlling for time-invariant plot attributes and time-varying site 677 attributes, the remaining variation in richness in a plot is "as if randomly assigned," 678 independently across time (note: a similar violation arises if P_{t-1} affects R_t , whereby ε_{t-1} is 679 necessarily correlated with R_t ; we address that possibility in the next section).

680 The theoretical and empirical literature on the causal effect of productivity on richness does 681 not have a clear conclusion: studies report productivity has zero effect on richness, a negative 682 effect on richness, and a humped-shaped effect. Nevertheless, there are some studies that report 683 detecting a negative effect of productivity on richness (e.g., (56)). To address this potential threat 684 to the internal validity of our estimated negative effect of richness on productivity, we take two 685 approaches: (a) we block a mechanism (M_t in Figure S4) through which productivity can affect 686 richness; and (b) we find a variable has no effect on productivity other than through its effect on 687 richness (Z_t in Figure S4) and use this instrumental variable to create an unbiased estimator of 688 the effect of richness on productivity.



689

690 **Figure S4. Reverse Causality in the Richness-Productivity Relationship.** In this causal graph, 691 richness in one period (R) has a causal effect on productivity (P) in the same period, and vice-692 versa. The variable M represents mechanisms that mediate these causal effects, which can be 693 different depending on direction of the causal arrow. Richness in the prior period (R_{t-1}) affects 694 richness in the current period (R_t). ε are time-varying factors that affect productivity, which, in 695 our estimation, we assume can be correlated across time. Z_t is often called an "instrumental 696 variable."

697

698 S6a.i. Blocked mechanism design

In regression analyses of a causal variable, it is well known that if one conditions on a mechanism variable, the estimated coefficient on the causal variable will no longer include the effect of the mechanism variable. In our grassland sample, if productivity were to negatively affect richness, we assume that this effect is, in part, mediated by shading; i.e., more productive plots generate greater shade, which in turn reduces richness (*56*). If the estimated negative relationship between richness and productivity in Figure 2 were an artifact of reverse causality mediated by shading, then putting our shading variable in Equation (S1) as a covariate would 706 block the effect of productivity on richness and the sign of the coefficient on richness (β) would 707 become positive (or small and statistically insignificant if the true relationship between richness 708 and productivity were zero). See Figure S5 below. If shading is not an important mechanism 709 through which productivity would affect richness in our sample, or if our measure of shading is a 710 poor measure of the shading mechanism, our mechanism-blocking design would fail to quantify 711 the potential threat of reverse causality. Indeed, productivity could alter biodiversity through 712 non-light pathways, such as soil resource use, but this effect of productivity on richness is 713 expected to, at least in part, be mediated by reductions in light from increased biomass that, in 714 turn, reduces richness in a plot. Thus, if reverse causality was a substantial threat to our 715 identification strategy for the effect of richness on productivity, we would expect the coefficient 716 of richness on productivity to shift towards more positive values.

As an estimate of shading, we measure the fraction of photosynthetically active radiation (e.g., light used by plants) that reaches the soil. This measure is calculated as the ratio of photosynthetically active radiation recorded below the plant canopy (ground level, mean of two readings) and that measured above the canopy. Measurements are carried out using a light meter (e.g. Ceptometer) at the same time and in the same 1m² sub-plots used for vegetation cover estimates. We have annual measures of ground-level light for 145 plots of our 151 plots (1011 of our 1231 observations).



724

- 528 but are not explicitly included in the graph.
- 729

Figure S5. Productivity affects species richness via shading as a mechanism. In this causal graph, in addition to changes in richness causing changes in productivity, changes in productivity change causes changes in richness via shading (s). Other mechanisms (M) may also be operative

730 First, we confirm that the estimated effect of richness on productivity does not change when 731 we use the subsample of 1,007 observations for which we have measures of shading. It does not: 732 a 10% increase in richness leads to an estimated 2.6% decrease in productivity, 95% CI [-4.4%, -733 0.8%] (see project page DOI/10.5281/zenodo.7675340). Next, we re-estimate Equation (S1) with 734 our shading variable included. The estimated negative effect of richness on productivity does not 735 change: a 10% increase in richness leads to an estimated 2.6% decrease in productivity, 95% CI 736 [-4.4%, -0.8%].² Said another way, if reverse causality was a substantial threat to our 737 identification strategy for the effect of richness on productivity, we would expect that, after 738 adding shading to Equation (S1), the coefficient of richness on productivity would become 739 substantially smaller in absolute value or, possibly, to become positive. Yet the estimate remains 740 unchanged.

741 S6a.ii. Instrumental variable design

As an alternative approach to assess the potential threat of reverse causality that makes *different* assumptions from the mechanism blocking analysis, we adopt another statistical approach that is common in economics and public health, but rare in ecology: an instrumental variable design (*21*, *22*, *57–59*).

746 We seek an attribute of the system that has a relationship with richness, but, after 747 conditioning on other attributes, has no relationship with productivity other than through its 748 relationship with richness. Such an attribute is illustrated by Z in Figures 1, S4 and S5. In 749 economics and biostatistics, Z is called an instrumental variable (IV) or a surrogate variable. An 750 example of a potential IV is randomization of planted richness by an experimenter. In field 751 experiments, randomization of richness helps isolate the causal effect of richness on productivity, 752 but only when the randomization affects productivity in a plot solely through its effect on 753 richness, an assumption called excludability (60) or the exclusion restriction (i.e., one must 754 assume there is no arrow going from Z directly to P).

In the absence of randomization, one must use theory and experience to identify a naturally
occurring IV. Each of the plots in our sample are unmanipulated plots that are embedded in

blocks of manipulated plots in the Nutrient Network. In other words, each unmanipulated plot in

² See the Github project page output for Table_MechBlocking_R_se.tex and Table_MechBlocking_R_ci.tex for more details.

758 our sample is surrounded by a set of plots with experimental nutrient additions (see (61)). These 759 manipulated experimental plots received randomized amounts of nutrient additions, which 760 subsequently affected the experimental plots' richness (62). We assume that the experimentally 761 manipulated richness in these plots can affect the richness in unmanipulated plots in the same 762 block through ecological dispersal channels but does not affect the productivity of these 763 unmanipulated plots except through the effect on the plots' richness (an assumption made more 764 plausible by the randomization of nutrients in the neighboring plots). If that assumption is 765 correct, we can use the average richness of an unmanipulated plot's neighboring manipulated 766 plots in the same block as an instrumental variable for richness in the unmanipulated plot.

767 This time-varying spillover effect from manipulated to unmanipulated plots is plausible if 768 either (a) the spatial pattern of nutrient manipulations in the experimental plots around each 769 unmanipulated plot varies across blocks, or (b) the spatial pattern of manipulated plots around 770 each unmanipulated plot varies across blocks (i.e., variation in how far apart plots in a block are 771 to each other or in their plot-level attributes that moderate the effect of nutrient addition n on 772 richness). We lack digital maps of the experimental designs for every site that we could use to 773 determine the exact distance between and spatial configuration of plots and empirically confirm 774 either of these assumptions. However, colleagues who manage the Nutrient Network believe 775 these assumptions are credible (Dr. Eric Seabloom, *personal communication*). Moreover, when 776 we regress unmanipulated plot richness on the average richness of the manipulated neighboring 777 plots in the block, we obtain a positive and statistically significant coefficient, which is 778 consistent with the posited spillover effect (Table S7, *column 2*).

779 **Table S7. Results from the Instrumental Variable Design**

	(1)	(2)
	Second Stage of 2SLS	First Stage of 2SLS
	(Outcome=Productivity)	(Outcome=Richness)
ln(Richness)	-0.24	
	(0.37)	
	[-0.96, 0.49]	
ln(Average		0.49
Neighboring Plots		(0.12)
Richness)		[0.26, 0.72]
Number of Plots	151	151
Number of Sites	43	43

Number of Observations	1212	1212
Montiel-Pflueger effective		
<i>F-statistic</i>		17.44

2SLS refers to a two-stage least squares estimator, with the results from the first stage (predicting richness) in column 2 and the results from the second stage (estimating effect of richness on productivity using the instrument from the first stage) in the column 1. The M-P effective F-statistics is used to test for a weak instrument (a test that is robust to heteroscedasticity, serial correlation, and clustering; (63)). The value of the M-P effective F-statimplies we can reject the null hypothesis of a weak instrument. Standard errors in parentheses (clustered at plot level) and 95% CI in brackets.

780 The excludability assumption implies that, after we condition on time-invariant plot attributes 781 and time-varying site attributes, the richness of a plot's manipulated neighbors affects the plot's 782 richness but has no effect on the plot's productivity other than via the effect on the plot's 783 richness. In other words, the drivers that cause the average richness of an unmanipulated plot's 784 neighbors to change over time only affect the unmanipulated plot's productivity through a 785 change in the unmanipulated plot's richness. The manipulated plots are randomly manipulated 786 (32) and these manipulations have been shown to affect species richness (62). Thus, some of the 787 changes in richness in neighboring manipulated plots are being driven by exogenous factors that 788 could plausibly be assumed to not affect unmanipulated plot productivity except through their 789 effects on the unmanipulated plots' richness. For this assumption to be valid, we must assume 790 that the nutrient additions that affect the neighboring plots' richness have no direct effect on an 791 unmanipulated plot's productivity other than via effects on an unmanipulated plot's richness 792 (e.g., rather than dispersal being the mechanism through which neighbor plot richness affects 793 own plot richness, it could be the nitrogen applications leaching through the ground). Colleagues 794 who manage the Nutrient Network believe this assumption is credible (Dr. Eric Seabloom, 795 personal communication).

796 In addition to the excludability assumption, we need two other assumptions to use this IV to 797 estimate a causal effect of richness on productivity: (1) first-stage non-zero effect of the IV; and 798 (2) monotonicity. The first-stage non-zero assumption of the IV design requires that there be a 799 correlation between neighbor's richness and own richness, on average - in other words, the 800 effect of neighbor's richness on own richness is not zero for all plots. We can verify this 801 assumption empirically (Table S8, column 1): after controlling for plot-level and site-level 802 confounders, the average richness of the neighboring plots has a positive association with own 803 plot richness. The monotonicity assumption implies that, for all plots, the relationship between a 804 plot's neighbor richness and its own richness can only be in one direction: it is either ≥ 0 or ≤ 0 . 805 In other words, we assume that we could not observe that, for some plots, higher neighbor

806 richness increases own richness, but for other plots, higher neighbor richness decreases own 807 richness. The monotonicity assumption is untestable. Yet given our ecological motivations for 808 using neighboring richness as an instrumental variable, we believe a non-negative, monotonicity 809 assumption is a valid approximation of the field reality, i.e., assuming that an increase in the 810 richness of surrounding manipulated plots can never decrease an unmanipulated plot's richness. 811 In comparison to our main design (Section S4), the IV design has two disadvantages. First, 812 because the IV design uses only variation in richness that comes from neighboring plot richness, 813 it will tend to have lower statistical power. Second, the IV design increases internal validity at 814 the potential expense of external validity. Rather than estimate the average effect on productivity 815 from any change in richness, we estimate the average effect for a subset of the changes in 816 richness. This subset is comprised of what are called "compliers" – plot-year observations for 817 which the richness value would have been different had the average richness in surrounding plots 818 been different. Given our instrumental-variable is multi-valued, there are many types of 819 compliers (e.g., plots that had 5 species in 2007 that would have had 6 species had their 820 neighboring plots had higher species richness). If the average causal effect of changes in richness 821 that come from changes neighboring plot richness differs from the average causal effect of 822 changes in richness that come from other attributes of the system, the generalizability of the 823 inferences from the IV design is more limited than in our main design. Another way to view the 824 more limited external validity of the IV design is that the IV design allows us to estimate the 825 average effect of richness on productivity for changes in a plot's richness that are induced by 826 changes in a neighbor's richness. If that average effect is not the same as the overall average 827 effect, the estimates from the IV design and the main design may differ, even if there is no 828 reverse causality or other forms of hidden bias in the main design.

In the IV design, we use two-stage least squares, linear, additive, fixed-effects estimator (*3*), which we implement with the "ivreghdfe" command in STATA v.16. The estimated effect of richness on productivity, as well as the first-stage estimates and F-test are reported in Table S8. To ensure reproducibility and use open-access software, we also estimate the IV equation in R using 'feols' in the fixest package (v. 0.8.2).

The estimate from the IV design implies a nearly identical estimate of the relationship between richness and productivity as we obtained from the main design: a 10% decrease in richness leads to a 2.4% increase in biomass. As expected, the effect is estimated imprecisely 837 (i.e., large confidence intervals). IV is a less efficient estimator (*12*), thus leading to predictably838 large confidence intervals.

839 Note on interference among plots: An unstated assumption in each design we implemented 840 in our study - an assumption found even in randomized controlled experimental designs - is "no 841 interference among units" (see (64)), which means that the potential productivity outcome in an 842 unmanipulated plot at a specific level of richness is unaffected by the richness levels of other 843 unmanipulated plots. The IV design may seem to imply that this assumption is violated, but the 844 IV design relies on manipulated plot richness affecting unmanipulated plot richness, which is not 845 interference in our study designs. In our designs, each unmanipulated plot is still assumed to 846 have one potential outcome per richness level. We believe that interference is absent in our 847 designs because, based on discussions with Nutrient Network coordinators, the unmanipulated 848 plots are sufficiently separated from each other within each site in order to not interfere with 849 each other. In other words, each unmanipulated plot's potential productivity outcomes under 850 different richness values only depends on its richness value, and not on the richness values of 851 other unmanipulated plots.

852 S6bi. Dynamic panel designs

In Figure S6, we present a more complicated causal graph than the DAG in Figure 1B. To make the new graph more compact, we do not specify whether variables are acting at the plot or site level, and we assume that all variables are measured in the current period, unless otherwise stated. In this new causal graph, we have productivity (P_t) and lagged productivity (P_{t-1}), we have common causes (U) of richness (R) and P_t and P_{t-1} , and common causes (I) of P_t and P_{t-1} . This causal graph is "unidentified" – in other words, there is no observational design that could estimate the effect of R on P without bias, unless we make more assumptions.



861 Figure S6. Reverse Causality in the Richness-Productivity Relationship in which Prior

- 862 **Productivity Affects Current Richness.** In this causal graph, richness in the current period (R_t)
- has a causal effect on productivity (P_t) in the same period, and productivity in the prior period
- 864 (P_{t-1}) has a causal effect on richness and productivity in the current period. The graph includes
- the two types of confounders that are the focus of the Main Design: time-invariant, plot-level
- 866 confounders (U_p) and time-varying, site-level confounders (U_{st}) . For simplicity, the graph does
- 867 not include the variable *I* from Figure 1 (i.e., factors that affect *P* but not *R*).







870 **Figure S7. Causal Graphs that Reflect the Main Design Estimation Strategy.** Our Main Design is valid for each causal graph in

871 this figure. The graph includes the two types of confounders that are the focus of the Main Design: time-invariant, plot-level

872 confounders (U_p) and time-varying, site-level confounders (U_{st}) . For simplicity, the graph does not include the variable *I* from Figure 1 873 (i.e., factors that affect *P* but not *R*). Our Main Design assumes the most important sources of confounding are from U_p and U_{st} and that a directed edge from P_{t-1} to R_t does not exist in combination with a directed edge from P_{t-1} to P_t (i.e., all panels in Figure S7 are allowed). Productivity can be serially correlated over time from unobserved variables, but there cannot be a causal effect of P_{t-1} on P_t that is mediated by R_t in this design.

879 If there were instead a directed edge from P_{t-1} to R, but either (i) no directed edge from U to

880 *R*, P_{t-1} and P_t (panel A of Figure S8; i.e., no unobserved common causes of richness and

productivity) or (ii) no directed edge from I to P_{t-1} and P_t (panel B of Fig. S8; i.e., no

unobserved, persistent causes of productivity), one could estimate the effect of *R* on *P* without

bias by conditioning on P_{t-1} ; in other words, via a lagged-dependent variable specification (in this

context we have both a form of reverse causation and a form of time-varying confounding).



885

Figure S8. Reverse Causality in the Richness-Productivity Relationship in which Prior Productivity Affects Current Richness. In these causal graphs, richness in the current period (R) has a causal effect on productivity (P) in the same period, and productivity in the prior period has a causal effect on richness and productivity in the current period. The graph includes the two types of confounders that are the focus of the Main Design: time-invariant, plot-level confounders (U_p) and time-varying, site-level confounders (U_{st}). For simplicity, the graph does not include the variable I from Figure 1 (i.e., factors that affect P but not R).

The causal processes implied by the causal graphs in Figure S7 and S8 are observationally indistinguishable. The data alone cannot tell us which estimation strategy, our Main Design or a lagged-dependent variable specification, is appropriate (unless one is willing to make untestable, parametric assumptions). In other words, productivity is temporally correlated across time – we can see that correlation in the data. The source of that temporal dependence could be unobserved persistent causes of productivity (often called "unobserved heterogeneity" in social science and

899 biostatistics). The design used to generate the main estimate in Fig. 2 controls for these causes 900 and thus eliminates any potential biases when such causes also are linked to richness. However, 901 the source of temporal dependence could be a direct link between productivity in one year and 902 productivity in the next year (e.g., via nutrient storage in roots). If that were the case, and if 903 biomass in one year also affected richness in the next year (often called "state dependence" in 904 social science and biostatistics), our main design may have bias. If, for example, productivity 905 was positively correlated across years and lagged productivity had a negative effect on current 906 richness, our estimated effect of richness on productivity using our main estimator would be too 907 negative. To eliminate that bias, we could condition on lagged productivity:

908
$$\ln LiveMass_{pst} = \beta \ln Richness_{pst} + \Omega \ln LiveMass_{ps(t-1)} + \mu_{st} + \varepsilon_{pst}$$
(S2)

909 In other words, if Fig. S8 is the correct interpretation of the systems we are studying, and we 910 use the estimator from Equation S1, the estimated β will be too large in the negative direction 911 (away from zero) compared to the true effect. If, however, Fig S7 is the correct interpretation of 912 the systems we are studying, and we use the lagged dependent variable estimator from Equation 913 S2, the estimated β will be too large in the positive direction (towards zero) compared to the true 914 effect. If the confounding process is a mix of the two, the two estimates bracket the true expected 915 causal response (3). We estimate Equation S2 using the "reghdfe" command in STATA (v.16) 916 and the 'feols' command in R using the fixest package (v. 0.8.2).

As expected, the estimated effect from Equation S2 is less negative (smaller in absolute
magnitude), providing an upper bound on the effect of species richness on biomass (Table S8).
Still, this estimate implies a negative effect of richness on productivity, with the lower 95%
confidence interval overlapping the estimate from the estimator in the Main Design (Equation
S1).

Finally, what if the true system were characterized by Fig. S6 and not approximated by either Fig. S7 or Fig. S8? Then we would have to make more assumptions to identify the effect of richness on productivity. For example, if we are willing to assume that the persistent causes of productivity across years (I) comprises autoregressive disturbances of order 1, the effect of R on P can be estimated using an autoregressive distributed lag equation of order 2 in autoregression and order 1 in distributed lags (65). Using the "reghdfe" command in STATA (v.16) and the 'feols' command in R using the fixest package (v.0.8.2), we estimate the following model:

929
$$\ln LiveMass_{pst} = \beta \ln Richness_{pst} + \beta' \ln Richness_{ps(t-1)} +$$

930
$$\Omega \ln LiveMass_{ps(t-1)} + \Omega' \ln LiveMass_{ps(t-2)} + \mu_{st} + \varepsilon_{pst}$$
(S3)

Table S8. Results from the Dynamic Panel Design. Results present coefficient estimates of a
1% increase in the ln richness on ln of productivity (measured as live biomass); clustered robust
standard errors are shown in the parentheses, and 95% confidence intervals in brackets. To see
95% confidence intervals as well, see the project page (TableS8_R_CI.tex).

n(SR)	-0.1329* (0.0789)	
	[-0.2876, 0.0218]	
$\ln LiveMass_{ps(t-1)}$	0.1374***	
	(0.0377)	
	[0.0634; 0.2115]	
Num. obs.	1063	
R^2 (full model)	0.83	

950 Signif. Codes: ***: 0.01, **: 0.05, *: 0.1

951 Robust Standard errors in parentheses (clustered at plot level).

When we estimate this equation S3, our estimated effect of richness on productivity is similar to our main estimate in Fig. 2: a 10% increase in richness leads to a 2.4% decrease in productivity, 95% CI [-4.4, -0.4]. Adding one more lag for richness and productivity (i.e., including ln $Richness_{ps(t-2)}$ and ln $LiveMass_{ps(t-3)}$) yields a similar estimate of -2.0% for ln *Richness_{pst}*. As a final estimator, we follow the suggestion of a peer reviewer and combine the Main Design estimator (Equation S1) and the Lagged Dependent Variable design (Equation S2). This

960 estimator is potentially biased for reasons that can be found in other publications (e.g., (3),

- 961 Sections 5.3-5.4). The intuition is that the combination of δ_{ps} from Equation (S1) and
- 962 $LiveMass_{ps(t-1)}$ from Equation (S2) can create a correlation between richness and the error
- 963 term in the model, which makes the estimator inconsistent (i.e., the estimator will not converge
- 964 in probability to the true value of our target parameter even as the number of plots in our panel

⁹⁵²
goes to infinity). This problem is not caused by an autocorrelated error process. The problem
arises even if the error process is *i.i.d.* If the error process is autocorrelated, the problem is even
more severe. Despite this potential bias, the estimated effect is similar to the estimates from the
other approaches: a 10% increase in richness, on average, decreases productivity by 2.8% [95%
CI: -4.6%, -1.1%].

970 <u>Note</u>: In Figures S6-S8, we do not include a graph in which $P_{t-1} \rightarrow R_t$, $R_t \rightarrow P_t$, and no other 971 direct edges exist. We exclude this graph because this pattern is not ecologically possible (i.e., a 972 case where prior productivity had an effect on current productivity and richness would be the 973 sole mediator or, equivalently, prior productivity has no effect on current productivity except 974 through its effect on current richness).

975 S6b.ii. Design sensitivity to unobserved, plot-level confounding variables

976 Imagine a set of unobserved, time-varying, plot-level confounding variables that are 977 negatively associated with richness and positively associated with productivity (so-called 978 "negative selection bias"). Were such a set of confounding variables to exist in our system, the 979 true average causal response in our sample could be closer to zero or positive. However, as we 980 will see below, using time-invariant plot attributes and time-varying site attributes, we can 981 explain about 90% of the variation in annual richness and productivity. Thus, there is little 982 variation left that could come from this unobserved variable and induce bias in our estimator. 983 Nevertheless, we test the sensitivity of our inferences to such a confounding variable.

984 Following the method introduced by Altonji et al. (66) and further developed by Oster (67), 985 we assume that the relationship between species richness and the omitted, unobserved 986 confounding variables can be characterized using information from the relationship between the 987 richness variable and the variables in Equation (S1). Oster in (67) assumes that observable 988 covariates contribute (approximately) proportionally to treatment and to the outcome (i.e., the 989 time-invariant plot attributes and the time-varying site attributes are as important in explaining 990 richness as they are in explaining productivity conditional on richness). If that assumption is a 991 reasonable approximation of the truth, which seems credible given our data (Table S7), we can 992 then explore the conditions under which the omitted variables could change our conclusions. 993 Said another way, this analysis answers the question, "How much correlation between the 994 unobserved variable and the richness/productivity variables would be sufficient to change our 995 conclusions?"

896 Recall that our preferred estimation Equation (S1) includes variables capturing time-invariant 997 plot attributes (δ_p) and time-varying site attributes (μ_{st}). That equation is most easily and 998 efficiently estimated by taking first differences or deviations from means, so that δ_p drops out 999 before estimating the coefficient β_1 . Thus, any goodness-of-fit measure, like R², would not 1000 include the role of the plot fixed effects in explaining variation in productivity.

1001 However, estimating Equation (S1) is equivalent to estimating the following specification:

1002
$$\ln (Biomass)_{pst} = \beta_0 + \beta_1 \ln (Richness)_{pst} + \mu_{st} + \sum_{ps} \alpha_{ps} d_{ps} + \epsilon_{pst}$$
(S4)

1003 where the variables are as before but rather than eliminate <u>time-invariant</u> plot attributes (δ_{ps}) by 1004 taking deviations from the mean (*Section S4*), we model them through a new set of variables: d_p 1005 is a set of plot dummy variables and α is a vector of coefficients of the plot dummy variables. 1006 We can estimate this equation using the 'areg' command in STATA (v.16).

1007 The estimated β_1 will be the same as the estimated β from our deviations-in-means 1008 estimation of Equation (S1), although the estimated standard errors will be larger using equation 1009 (3) (i.e., estimating Equation (S4) is less efficient than estimating Equation (S1) via deviations in 1010 means). These differences in estimated standard errors come from clustering the standard error 1011 estimation at the plot-level. Clustering standard errors in the 'areg' procedure adjusts the degrees 1012 of freedom by the number of fixed effects removed in the within-group transformation. In 1013 contrast, the reghdfe (or xtreg fe) procedure does not make such an adjustment and thus reports 1014 smaller cluster-robust standard errors. Thus, to estimate the effect of richness on productivity, the 1015 xtreg fe or reghdfe estimation procedure is preferred. However, for our sensitivity analysis, using 1016 the 'areg' procedure is preferred because it allows us to use the time-invariant, plot-level 1017 characteristics as control variables, rather than treat them as nuisance parameters that can be 1018 eliminated via first-differencing or taking deviations from means. Table S9 (column 1) reports 1019 the results from the areg regression with clustered standard errors at the plot level. As expected, 1020 the 'areg'-estimated coefficient on species richness is the same as the estimate in the main text, 1021 but with a larger standard error.

1022To explore the sensitivity of our design to an unobserved confounder, or set of confounders,1023Oster (67) shows that one must first make an assumption about the R-squared value from a1024hypothetical regression of productivity on the unobserved confounder and the variables in1025Equation (S4). If this value, called R_{max} , were set equal to one, it would be equivalent to saying

- 1026 that the variation in annual productivity can be fully explained by the hypothetical regression
- 1027 model. In other words, the unobserved confounder explains all the unexplained variation in
- 1028 regression specification (S4), an assumption that implies there is no measurement error in the
- 1029 live biomass measure of annual plot productivity. The R-squared from the regression
- 1030 specification (S4) is 0.87 (Table S9, column1). Thus, an assumption that $R_{max} = 1$ implies that the
- 1031 unobserved confounder explains all remaining variation in annual productivity, which is an
- 1032 implausibly powerful predictor of productivity (i.e., an implausibly strong confounder).

	(1)	(2)
	Outcome Equation:	Selection Equation:
	Productivity	Species Richness
Species Richness	-0.24	
	(0.09)	
	[-0.42, -0.06]	
R-Squared	0.87	0.91
Number of Plots	151	151
Number of Sites	43	43
Number of Observations	1231	1231

1033 Table S9. Regression results to support sensitivity analysis to hidden bias.

The first column replicates our main result from column 1, Table S2 but uses a dummy variable procedure to control for plot-level fixed effects rather than a deviations-from-means procedure (and thus standard error estimates are slightly larger in this table). The second column regresses plot-level richness on plot dummy variables and the year-by-site variables. The sample size is larger in second column because fewer plots have missing richness values than have missing productivity values. Robust standard errors in parentheses (clustered at plot level) and 95% CI in brackets.

1035 We next specify the magnitude of the degree of selection on unobservable variables relative 1036 to the selection on observable variables. This parameter π (called δ in Oster's article) yields a 1037 measure of our design's sensitivity to hidden bias: how would our results change were there an 1038 unobserved confounder (or set of confounders) that is correlated with richness and productivity? 1039 The unobserved confounder is assumed to be uncorrelated with the other variables in our equation; it is most easily envisioned as a time-varying confounder that is orthogonal to our site-1040 1041 by-year variables and our time-invariant plot variables. To help determine a plausible value for π , 1042 we use 'areg' to estimate the relationship between the richness and the observable covariates in 1043 Equation (S4) using the following specification:

1044
$$\ln (Richness)_{pst} = \beta_0 + \mu_{st} + \sum_{ps} \alpha_{ps} d_{ps} + \epsilon_{pst}$$
 (S5)

¹⁰³⁴

1045 Results from this regression are reported in Table S9. The R² from this regression (Table S9, 1046 column 2) implies that the variables in our regression specification already explain about 91% of 1047 the variation in annual plot richness. Setting the value of π equal to -0.10 would be equivalent to 1048 assuming that the unobserved confounder explains all the remaining variation in richness (i.e., 1049 the unobserved confounder explains about 9% of the variation in annual plot richness).

1050 Setting $\pi = -0.10$ and $R_{max} = 1$ would mimic a powerful potential unobserved confounder in 1051 our design: a confounder that is so strongly correlated with productivity and richness that, were 1052 we able to observe it (along with the other variables in the equation), we could predict with near 1053 certainty which of two plots would have higher productivity and which would have higher 1054 richness. Estimating the effect of richness on productivity with those implausible parameter 1055 values yields an upper bound on the impact of richness on productivity. For completeness, we 1056 also calculate a lower bound on the impact by setting $\pi = 0.10$ (i.e., positive selection bias). 1057 To estimate these bounds, we use the 'psacalc' package in Stata (v.16) as a post-estimation 1058 function after estimating the effect of richness on productivity using Stata's 'areg' command. 1059 The estimated upper bound is still negative: a 10 % increase in richness implies a 2.0% decrease 1060 in productivity. In other words, in the presence of an unobserved confounder that is negatively 1061 associated with richness and positively associated with productivity relationship (thus creating 1062 some spurious correlation between richness and productivity), we would still infer that there is a 1063 negative relationship between richness and productivity. To infer a positive relationship between 1064 the two variables would require an infeasible value for π : it requires $\pi > 1$, which implies the 1065 confounder would have to be more influential in explaining variation of productivity than the 1066 plot-level, time-invariant attributes and the site-level, time-varying attributes. We also consider 1067 an unobservable confounder that is masking some of the negative effect of richness on 1068 productivity. If the unobserved confounder was positively associated with both richness and 1069 productivity, a 10% increase in richness would decrease productivity by an estimated 3.0%. 1070 Overall, the sensitivity analysis implies that our estimated effect of richness on productivity 1071 is not sensitive to the presence of a time-varying confounder.

- 1073
- 1074

1075 S7. Comparison of Main Design to Common Designs in Ecology

Here, we refer to multi-level modeling, hierarchical modeling, random effects modeling,
mixed effects modeling, or variance components modeling as a Common Design in Ecology with
which we compare our Main Design. Ecologists who are familiar with multi-level modeling
may wonder why, given our data are comprised of plots nested within sites and annual
observations nested within plots, we do not use this multi-level modeling approach as our main
design. A full explanation is beyond the scope of this SI, but the main reason has two parts:

1082 (a) Without more variable transformations, the multi-level modeling approach does not 1083 easily lend itself to controlling for as many unobservable sources of confounding as can 1084 be done in our linear, additive, fixed-effects panel data estimator. Note that our use of the 1085 term 'fixed effects' corresponds to the use of the term in econometrics (3, 4) – and the 1086 meaning differs from the use in multi-level modeling uses in Ecology; see Glossary S1 1087 for more details. The typical multi-level model assumes "selection on observables" (see 1088 Glossary S1). In other words, the model requires stronger assumptions to infer causality 1089 from the correlation between richness and productivity; it assumes that there is no 1090 correlation between unobserved components of the error term at the site and plot levels 1091 and the species richness variable (17, 44). This assumption is easily violated in this 1092 dataset. In other words, the multi-level model assumes that the propensity of plots to 1093 experience change in their species richness is not determined by variables that also affect 1094 productivity and are not explicitly in our model (i.e., not determined by variables in the 1095 error term). In contrast, in our panel data design, the time-invariant individual and site 1096 attributes are no longer part of the error term. Thus, we do not have to make that strong 1097 "selection on observables" assumption made by traditional multi-level models.

(b) Clustering our standard errors at the plot provides the same benefits that multi-level
 modeling does when estimating the variance-covariance matrix in the presence of intra site correlations among plots (44).

1101 The disadvantage of our design (i.e., Equation S1) is that, in the process of eliminating the 1102 role of confounding variables, we also eliminated variables that affect productivity but do not 1103 affect richness (I in Figure 1B) – i.e., both the "bad" between-plot (confounders) and "good" 1104 between-plot variation (predictors) are eliminated, thus reducing statistical power (increasing standard error estimates) and limiting the scope of some of the analysis that can be performed (68). For example, we cannot, in our study, estimate the associations between productivity and time-invariant attributes like soils or climate. Some other research aims, like estimating how the effect of richness on productivity varies by site (versus site-level moderators), are easier to accomplish in multi-level modeling designs.

1110 In our study, however, we are focused on estimating the causal effect of biodiversity on productivity. Thus, in our study context, the traditional multi-level modeling approach – which 1111 1112 aims to control for observable confounding variables by including them in the model - has no 1113 advantages over our approach, and a serious disadvantage in its inability to control for 1114 unobservable sources of confounding variation without data transformation that mimics what 1115 happens in our estimation of Equation (S3) (see last paragraph in this section for details). 1116 Indeed, not controlling for unobservable sources of confounding variation leads to bias in 1117 estimates.

1118 Nevertheless, to demonstrate how much design matters in drawing inferences from1119 observational data, we estimate a traditional random-effects equation by using a GLS estimator

1120
$$\ln LiveMass_{pst} = \beta \ln Richness_{pst} + \sum \gamma X_p + \sum \theta X_s + \vartheta_t + \varphi_p + \rho_s + \varepsilon_{pst}$$
 (S6)

1121 where φ_p and ρ_s are parts of the error term and X_p and X_s are observable attributes of plot and site, respectively. This equation represents the "Common Design in Ecology" that we report 1122 1123 in the main text. Unlike Equation (S1), this equation uses both the within-plot variation in 1124 species richness and the between-plot variation in richness to estimate the effect of richness on 1125 productivity (i.e., the estimator produces a matrix-weighted average of the between-plot and 1126 within-plot estimates). The between-plot variation in species richness can be caused by many 1127 attributes that may also be correlated with productivity. We follow a traditional approach in 1128 ecology and attempt to control for (block) the confounding effects of these attributes by 1129 measuring them and adding them to the equation $(X_p \text{ and } X_s)$. In the words of the multi-level 1130 modeling community, the γ and the θ are "fixed effects" and the φ_p and ρ_s are "random effects." 1131 The "fixed effects" are directly estimated whereas the "random effects" are not, but rather 1132 summarized according to their estimated variances and covariances.

Because our data set is large, we can control for over 60 observable confounders, including variables for soil attributes, habitat types, historical management categories, weather, year, country, and elevation (see Table S10). In other words, we can control for a wide range of variables that may affect both richness and productivity. Along with species richness, these variables explain 57% of the overall (spatial and temporal) variation in productivity (live biomass) and 95% of between-plot (spatial) variation in productivity.

Using this design, the estimated relationship between richness and productivity is much more in line with the conventional wisdom in the ecology literature: a 10% increase in plot richness increases plot productivity, on average, by 3.8%, 95% CI [2.0%, 5.6%]. In other words, when we do not leverage the spatial and temporal variation to eliminate the effects of unobserved plot and site attributes (i.e., our Main Design), we draw the opposite conclusion about the relationship between richness and biomass.

1145 The problem with this *Common Design in Ecology* is too few control variables, not too 1146 many. Failing to control for a confounding variable leads to omitted variables bias (reviewed in 1147 (69)), yet it is impossible to know or measure all confounding variables in a complex ecosystem 1148 system. Thus, the problem cannot be solved by adopting a "model selection" procedure to select, 1149 based on some measure of prediction error, a subset of the 60 variables (e.g., forward-selection 1150 or backward-selection procedures). Including control variables that have no correlation with the 1151 outcome will indeed add noise to the estimation procedure and unnecessarily reduce the 1152 precision of the estimated effect of richness on productivity. But precision is not a problem in 1153 our study because of our large sample size. Potential bias is the problem. Thus, the choice of 1154 control variables must not be driven by statistical benchmarks but rather by theory and field 1155 experience about which ecosystem features may affect both richness and productivity. One can 1156 only justify eliminating control variables from the model if one believes that the remaining 1157 covariates eliminate the correlation between richness and the model error term, which can only 1158 be justified based on theory or field experience.

A multi-level structural equation model (SEM) would be more flexible than the multi-level regression we implement in Equation (S6), but if we used the same observable variables, it would not affect the coefficient estimates. Use of a SEM would affect the standard error estimates, which is not relevant to make our point related to how our inference on the sign of the estimated effect of richness on productivity switches across designs. SEMs are also useful for

- 1164 estimating mediator effects (e.g., by specifying direct and indirect paths), estimating the
- 1165 relationship between observed and latent variables, and developing predictive models of
- 1166 productivity, but none of those aims are relevant for our study design.

Attribute	Covariate Variables
Country variables	Australia (AU), Canada (CA), Switzerland (CH), Germany (DE), Tanzania (TZ), United Kingdom (UK), United States (US)
Habitat variables	Alpine grassland, Annual grassland, Desert grassland, Mesic grassland, Montane grassland, Old field, Pasture, Savanna, Semiarid Grassland, Shortgrass prairie, Shrub steppe, Tallgrass prairie
Observation year variables	1 st , 2 nd , 3 rd , 4 th , 5 th , 6 th , 7 th , 8 th , 9 th , 10 th , 11 th
(Year in site's panel data set)	
Historical site management variables	Active management (otherwise wild), Active managed burning regime, Regularly grazed by herbivores, Restored
Topographical variables	Elevation (meters)
Weather variables*	Temperature Seasonality (standard deviation *100), Max Temperature of Warmest Month, Min Temperature of Coldest Month, Mean Temperature of Wettest Quarter, Mean Temperature of Driest Quarter, Mean Temperature of Warmest Quarter, Mean Temperature of Coldest Quarter
Soil physical property variables	Soil Percent Sand, Soil Percent Silt, Soil Percent Clay
Soil fertility variables	Soil Percent Carbon by Mass, Soil Percent Nitrogen by Mass, Soil Phosphorus by Mass (ppm), Soil Potassium by Mass (ppm), Soil Calcium by Mass (ppm), Soil Magnesium by Mass (ppm), Soil Sulfur by Mass (ppm), Soil Sodium by Mass (ppm), Soil Zinc by Mass (ppm), Soil Manganese by Mass (ppm), Soil Iron by Mass (ppm), Soil Copper by Mass (ppm), Soil Born by Mass (ppm), pH

Table S10. Control Variables in Common Design in Ecology 1167

1168 * The weather variable values in NutNet are site-level averages over time (i.e., they are time-invariant). Thus, 1169 controlling for temperature variables listed in Table S10 also controls for the precipitation variables,

1170

evapotranspiration variables, and other weather variables. In other words, the temperature variables serve as site-1171 level indicator variables and the estimated effect would be the same if, for example, we used precipitation variables

- 1172 instead of temperature variables.
- 1173 The positive estimated effect from Equation (S6) is not driven by having to drop sites that did
- 1174 not measure all the covariates (the sites in France, Portugal, and South Africa did not collect the
- 1175 soil data). If we use only the 675 observations from the multi-level modeling in our Main
- 1176 Design, we still obtain a negative estimated effect of richness on productivity, albeit less
- 1177 precisely estimated because of the smaller sample size: a 10% increase in plot richness decreases

plot productivity, on average, by 3.11%, i.e., the estimated effect is - 3.11%, with 95%

- 1179 confidence interval of [-6.31%, 0.09%]. Thus, the contrast between the Main Design and the
- 1180 Common Design in Ecology is not affected by the change in the sample composition.

1181 The key issue that we highlight is one of "design" – not of methods (i.e., type of estimation 1182 procedure). In principle, one could accomplish the same objectives of our design within a multi-1183 level or SEM framework by using a group-mean centering transformation of the data (i.e., 1184 within-plot centering of time-varying richness; (68, 70)). The key innovation in our *design* is to 1185 leverage the panel data to control for a wide range of time-invariant plot attributes and time-1186 varying site attributes -- a wider range of confounders than previous studies have addressed. 1187 Whether that leverage is exploited in a single regression equation or in a system of regression 1188 equations matters little for the estimation of the effect of plot richness on plot productivity.

1189

1190 S8. Supplementary Methods: Heterogeneous Effects of Rare, Non-rare, and

1191 Non-native Species on Productivity

To shed light on the reasons why an increase in species richness reduces productivity, on average, we decompose species richness into groups of species. Using the plot-level data from the Nutrient Network (see section S8a), we first decompose overall species richness into native versus non-native species and rare versus non-rare species (by "non-rare", we mean dominant and subordinate species). We then create groups for four categories of species that combine rarity and non-native status: (1) rare and native, (2) rare and non-native, (3) non-rare and nonnative, and finally (4) non-rare and native.

1199 As in most ecosystems worldwide (71), and consistent with theory (72), the rank abundance 1200 curves of species from our sites imply that most species in these ecosystems are rare (Figure 1201 S10). Thus, not surprisingly, higher species richness at our 43 sites is associated with, on 1202 average, higher numbers of rare species (Fig S11(A)). Moreover, as in many natural ecosystems 1203 in the Anthropocene (73), higher species richness at our 43 sites is also associated with, on 1204 average, higher numbers of non-native species (Fig S11(B)). Thus, higher numbers of species 1205 tend to be associated with more rare and non-native species in a place, and these species may 1206 have different effects on productivity than native non-rare species (e.g. 60, 61). Our analysis

1207 contributes to our understanding of how the effect of species richness on productivity depends on1208 the characteristics of the biodiversity changing.

- 1209 All code to reproduce the data processing in section S8 can be found at the project page. This 1210 code was checked by 3 additional people beyond the lead author.
- 1211





1213 Figure S10. Most species are rare in these grassland ecosystems. Rank abundance curves

1214 (RAC) for each Nutrient Network site in our analysis shown in Table S1. These RACs are for the

1215 pre-treatment year, which we use to define species as rare or non-rare. Here, species abundance 1216 is based on relative live cover at the site-level.



(B) Higher SR is positively associated with non-native SR



1218

Figure S11. Greater diversity, in terms of more numbers of species, is associated with more
(A) native rare species and (B) non-native species, on average. We plot the association

between overall species richness (SR) and counts of native rare species and of non-native species, per plot and year.

1223 S8a. Definitions and Measurement of groups in Figure 5

To classify species by rarity and origin, we use plot-level data from the entire site (i.e., not just our 151 unmanipulated plots). The rarity designation is based on measures of relative abundance at the site. To ensure that our relative abundance measures are unaffected by the experimental manipulations at the site, we use data only from the pre-treatment years. Species absent during the first year were treated differently (see below). Note that because we classify the species in the groups based on pre-treatment year data, the

1230 site "saline.us" is excluded from this analysis because the site does not have pre-treatment data.

Dropping the 24 observations from the site "saline.us" does not change our estimates in Figure 2 and 3.

1233 Native versus Non-native Species

1234 In the Nutrient Network, species origin was determined by the site coordinators and

- designated by one of three categories: "native", "INT" (i.e., non-native), or "unknown origin"
- 1236 (see "ProcessNutNet_coverData_FINAL_public.R" code on our project GitHub release
- 1237 DOI/10.5281/zenodo.7675340 for more details). We compute the number of species that are
- 1238 classified as native and non-native in each plot in each year and then construct species richness

1239 variables for each species type in each plot in each year. We drop the species of "unknown

1240 origin" present in the pre-treatment year in our main analysis, but to consider how the

1241 uncertainty about the origin of some species in the data could affect our conclusions, we

1242 performed a bounding approach where we re-estimate the effects by first assuming all unknown

1243 origin species are native and then assuming they are all non-native (see *S8ci. Sensitivity*

1244 Analyses: Species with unknown origin).

1245

1246 Classifying Rare versus Non-Rare Species

1247 We assign the labels of "rare" and "non-rare" to species in multiple ways, using definitions 1248 based on two metrics for relative abundance of a species at a site: the relative cover of each 1249 species at a site and the relative frequency of each species at a site. To calculate relative cover 1250 and relative frequency, we use live cover only in the pre-treatment year. Next, we classify 1251 species as "rare" based on their cover and frequency relative to other species at each site (see 1252 below). We then compute the number of species that are classified as rare or non-rare in each 1253 plot in each year and then construct species richness variables for each species type in each plot 1254 in each year. Below, we describe each step in this procedure.

In the main text, we define "rare" species using the relative cover metric. We use relative cover as our metric for abundance, because we believe better captures the range of mechanisms through which rare species may decrease productivity, including taking space formerly occupied by more productive species. In *Section S8c.ii*, we report the results using the relative frequency metric. In *Section S8c.iii*, we also test the sensitivity of our conclusions to different cutoff values for assigning a species to the "rare" and "non-rare" categories.

1261 *Computing the relative cover metric*

For each species present in the pre-treatment year, we computed the relative cover as the sum of the plot cover of the species for all plots at the site, divided by the total cover in all plots at that site. Note that some plots exceed 100% cover, whereas other plots are <100%; thus, we standardize this metric by dividing the sum of the maximum cover of a species in each plot at a site by the total live cover in all plots at a site.

1267 Classifying species as rare or non-rare using relative cover

Using the relative cover metric, we classify species into three categories: dominant,
subordinate, and rare. The categorization of rare, subordinate, and dominant species is based on
the quantiles of the species' relative cover data for each site, created using the 'quantile' function
in R.

1272 The first classification, presented in the main text Figure 5 and Table S10, labeled species at a 1273 site with relative cover in the lowest 60% of the site-level distribution (0.6 quantile) to be rare 1274 and species in the top 95% of the distribution (0.95 quantile) to be dominant. The species with 1275 relative cover in between these two cut-off values were labeled subordinate. These cut-off values 1276 lead to a maximum of 1 to 4 species dominant species per site, consistent with (76), and a 1277 median of 2 dominant species (average of 2.3). This cut-off leads to an average of 21.7 rare 1278 species and median of 20 rare species per site and an average of 12.73 and median of 12 1279 subordinate species per site. To maintain objectivity in the analysis, the person who 1280 recommended these cut-off values (Kaitlin Kimmel) was not the person running the estimation 1281 analyses. To assess whether the cut-off values generated a sensible classification of species – 1282 particularly with regard to differentiating dominant from rare species – the person who 1283 recommended the cut-off values (co-author Dr. Kaitlin Kimmel) checked which species were 1284 labeled "dominant" at two sites about which she had extensive knowledge (cdcr.us and knz.us). 1285 She confirmed that the three species that were labeled dominant at each site were indeed what 1286 experts would label as the dominant species.

Given that few species are dominant at each site and, by definition, these species do not exit and enter plots with great frequency, we combine the numbers of dominant and subordinate species into the non-rare species richness variable. We test the sensitivity of our results to the choice of cut off values for classifying species as rare or non-rare (Section S8c.iii).

Several species were not observed in the first year of data collection at a site, implying that
those species had a relative cover and relative frequency of 0 in that pre-treatment year.
However, rather than assume these species are rare or non-rare, we classified these species

separately (as "NA species") and controlled for them in our analyses (see Section *S8c.i*).

1295 Once we classify species are "rare", "non-rare" and "NA" based on pre-treatment data, we 1296 then count the number of species in each combined category for each plot and year.

1297

1298 S8b. Statistical Analyses

Combining the classification of species by origin and the classification of species by rarity, we can count the number of Non-Rare Native species, Rare Native species, Non-Rare Non-Native Species, Rare Non-Native species, and species classified as NA for each plot and year. We then substitute the richness variables of each species category for the overall "species richness" variable used in our main Equation (S1). In other words, we substitute the five categories of richness for the single "richness" variable of Equation (S1):

 $1305 \quad \ln(LiveMass_{pst}) = \beta_{DN} ihs(NonRareNative_{pst}) + \beta_{RN} ihs(RareNative_{pst}) + \beta_{DNN} ihs(NonRareNative_{pst}) + \beta_{RNN} ihs(RareNonNative_{pst}) + \beta_{NA} ihs(NAspp_{pst}) + 1307 \quad \delta_{ps} + \mu_{st} + \varepsilon_{pst}$ (S7)

The species richness variables are transformed with an inverse hyperbolic sine transformation rather than a natural logarithm transformation. The inverse hyperbolic sine transformation is analogous to a log-transformation but can be used when there are many 0 observations (*52*). Given the inverse hyperbolic sine transformation of the richness variables, the estimated effects cannot be interpreted as elasticities without further manipulation, but their signs and relative magnitudes can be compared to each other.

1314 S8c. Comparing the effect of species richness per group on productivity

1315 In Figure 5 and Table S11, we present the estimated effects of species richness on 1316 productivity, conditional on species type, using relative cover to classify species as non-rare or 1317 rare. The estimates imply that an increase in species richness has a positive effect on productivity 1318 when the increase is coming from non-rare, native species or from rare, non-native species. But 1319 for rare, native species, as well as for non-rare, non-native species, the estimated effects are 1320 negative. The estimated effect is largest in absolute value for rare, native species and non-rare, 1321 non-native species. We reject the null hypotheses that changes in the richness of these groups of 1322 species have an equivalent effect on live mass (*Chisq* = 9.8205, *Pr*(>*Chisq*) = 0.02016).

Controlling for the number of species labeled "NA" that enter the plots after the pre-treatment year (which we call "NA species richness") does not change the inferences drawn (Table S11 – i.e., compare column *1 versus* column 2). In the main text Figure 5, we present the conservative model that controls for the NA species richness.

1328 Table S11. The effect of species richness (SR) on biomass production conditional on species

1329 **type,** using species relative cover to determine rare versus non-rare species. We estimate

Equation (S1) with species richness disaggregated into the numbers of non-rare native species,

rare native species, non-rare non-native species, and rare non-native species. In column 1, we

1332 controlled for the number of species not found in the pre-treatment year at site; in column 2, we1333 perform a sensitivity analysis, dropping that species count. All estimated effects of each category

- 1334 of species richness (SR) are on a log-inverse-hyperbolic-sine scale. Given the inverse hyperbolic
- 1335 sine (his) transformation of the species richness variables, the estimated effects cannot be
- 1336 interpreted as elasticities without further manipulation, but their signs and relative magnitudes

1337 can be compared to each other. Clustered-robust standard errors are used and clustered at the plot 1338 level. To see 95% confidence intervals as well, see the project page (output/TableS11 R CI.tex).

	Main Text	Dropping counts of NA species
Non-rare, Native SR	0.0488	0.0507
	(0.0721)	(0.0723)
Non-Rare, Non-Native SR	-0.1721***	-0.1794***
	(0.0649)	(0.0651)
Rare, Non-Native SR	0.0397	0.030
	(0.0711)	(0.0701)
Rare. Native SR	-0.1473***	-0.1429***
	(0.0459)	(0.0466)
SR of NA species	-0.0901**	
	(0.0430)	
Num. obs.	1,175	1,175
Num. plots	146	146
R^2 (full model)	0.79	0.79

1357 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

1358 Robust Standard errors in parentheses (clustered at plot level).

1359

1360 S8ci. Sensitivity Analyses for species with unknown origins

Several species were classified as of "unknown origin." The analyses presented in Table S11 and S13 and Figure 5 omit species of unknown origin from the groups. To test the sensitive of our results to this uncertainty, we bound the estimated effects by considering two possible extreme scenarios: 1) all species of unknown origin are native and 2) all species of unknown origin are non-native. Thus, we revise the species groups in Equation (S7) and re-rerun the analyses with these two sets of models to establish bounds. The signs and magnitudes of the estimated effects in Tables S12 and S13 are similar to those in Table S11. 1368 Table S12. Sensitivity Analyses: treating species of unknown origin as native. We estimate 1369 Equation (S1) with species richness disaggregated into the numbers of non-rare native species, 1370 rare native species, non-rare non-native species, and rare non-native species. Here we treat 1371 species of unknown origin as native (rare or non-rare). In column 1, we controlled for the number of species not found in the pre-treatment year at site; in column 2, we perform a 1372 1373 sensitivity analysis, dropping that species count. All estimated effects of each category of species 1374 richness (SR) are on a log-inverse-hyperbolic-sine scale. Given the inverse hyperbolic sine 1375 transformation of the species richness variables, the estimated effects cannot be interpreted as 1376 elasticities without further manipulation, but their signs and relative magnitudes can be 1377 compared to each other. Clustered-robust standard errors are used and clustered at the plot level. 1378 To see 95% confidence intervals as well, see the project page (output/TableS12 R CI.tex).

ng for NA species	Dropping counts of NA species
0.0848	0.0878
(0.0712) -0.1717***	(0.0719) -0.1796***
(0.0640)	(0.0646)
0.0367	0.0263
(0.0704)	(0.0693)
-0.1399***	-0.1333***
(0.0458)	(0.0470)
-0.0914**	
(0.0432)	
1,175	1,175
146	146
0.79	0.79
	ng for NA species 0.0848 (0.0712) -0.1717*** (0.0640) 0.0367 (0.0704) -0.1399*** (0.0458) -0.0914** (0.0432) 1,175 146 0.79

1397 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

1398 Robust Standard errors in parentheses (clustered at plot level).

1399 S8cii. Sensitivity analyses using relative frequency as a metric for rarity

1400 We next check the robustness of the results to our choice of metric for determining rarity,

1401 comparing our results in Figure 5 to estimates when defining rarity based on relative frequency,

1402 rather than relative abundance. We calculate relative frequency as the number of plots that a

species occurs in a year divided by the total number of plots at a site in each year.

1404 The inferences about the effects of different species groups on productivity are similar using

1405 frequency instead of cover (Table S14). The Zenodo release of project page also includes results

1406 for relative frequency dealing with unknown species origin as above in section *S8c.i.* Results are

similar; for more details see the project page (/output/TableS14Sensitivity_R_se.tex).

Table S13. Sensitivity Analyses: treating species of unknown origin as non-native. We estimate Equation (S1) with species richness disaggregated into the numbers of non-rare native species, rare native species, non-rare non-native species, and rare non-native species. Here we treat species of unknown origin as non-native (rare or non-rare). In *column 1*, we controlled for the number of analysis not found in the num treatment user at sites in *column 2*, we perform a

1412 the number of species not found in the pre-treatment year at site; in *column 2*, we perform a 1413 sensitivity analysis, dropping that species count. All estimated effects of each category of species

richness (SR) are on a log-inverse-hyperbolic-sine scale. Given the inverse hyperbolic sine

1415 transformation of the species richness variables, the estimated effects cannot be interpreted as

1416 elasticities without further manipulation, but their signs and relative magnitudes can be

1417 compared to each other. Clustered-robust standard errors are used and clustered at the plot level.
1418 To see 95% confidence intervals as well, see the project page (output/TableS13 R CI.tex).

Controlling	for NA species	Dropping counts of NA species
Non-rare, Native SR	0.0451	0.0458
	(0.0725)	(0.0726)
Non-Rare, Non-Native + unknown SR	-0.1362**	-0.1382**
	(0.0669)	(0.0666)
Rare, Non-Native + unknown SR	0.0512	0.0514
	(0.0617)	(0.0616)
Rare. Native SR	-0.1500***	-0.1467***
	(0.0456)	(0.0462)
SR of NA species	-0.0916**	
-	(0.0432)	
Num. obs.	1,175	1,175
Num. Plots	146	146
R^2 (full model)	0.79	0.79

1437 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

1438 Robust Standard errors in parentheses (clustered at plot level).

1439

1440

1441

1442

1443

1445 Table S14. Sensitivity Analysis using Relative Frequency. We determine species groups and 1446 the effect of species richness (SR) on biomass production conditional on species type, using 1447 species relative frequency to determine rare versus non-rare species. We estimate Equation (S1) 1448 with species richness disaggregated into the numbers of non-rare native species, rare native 1449 species, non-rare non-native species, and rare non-native species. In *column 1*, we controlled for 1450 the number of species not found in the pre-treatment year at site; in column 2, we perform a 1451 sensitivity analysis, dropping that species count. All estimated effects of each category of species 1452 richness (SR) are on a log-inverse-hyperbolic-sine scale. Given the inverse hyperbolic sine transformation of the species richness variables, the estimated effects cannot be interpreted as 1453 1454 elasticities without further manipulation, but their signs and relative magnitudes can be 1455 compared to each other. Clustered-robust standard errors are used and clustered at the plot level. To see 95% confidence intervals as well, see the project page (output/TableS14_R_CI.tex). Also, 1456 please find other sensitivity analyses as done in Table S12 and S13 but using relative frequency 1457 on the project page (analyses_fig5_smsection8.R). The conclusions remain unchanged. 1458 1 1 5 0

Controlling for NA species	Dropping counts of NA species
0.0411	0.0505
(0.0804)	(0.0825)
-0.1739**	-0.1844***
(0.0688)	(0.0685)
-0.0126	-0.0183
(0.0606)	(0.0592)
-0.0936*	-0.0914*
(0.0500)	(0.0500)
-0.0834*	
(0.0420)	
1,175	1,175
146	146
0.78	0.78
	Controlling for NA species 0.0411 (0.0804) -0.1739** (0.0688) -0.0126 (0.0606) -0.0936* (0.0500) -0.0834* (0.0420) 1,175 146 0.78

1477 Signif. Codes: ***: 0.01, **: 0.05, *:0.1

1478 Robust Standard errors in parentheses (clustered at plot level).

1479 S8ciii. Sensitivity Analyses using different cut-offs for rare versus non-rare categories

To assess the sensitivity of the results to the classification criteria for rare and non-rare species, we use two additional cut-offs. Cut-off 2 labels species at a site with relative frequency in the lowest 70% of the site-level distribution (0.7 quantile) to be rare and species the top 95% of the distribution (0.95 quantile) to be dominant. The species with relative cover in between these two cut-off values were labeled subordinate. Cut-off 3 labels species at a site with relative frequency in the lowest 50% of the site-level distribution (0.5 quantile) to be rare and species the top 95% of the distribution (0.95 quantile) to be dominant. The species with relative cover in the lowest 50% of the site-level distribution (0.5 quantile) to be rare and species the between these two cut-off values were labeled subordinate. As in the analysis for Table S11, the
subordinate and dominant species were grouped together in the "non-rare" category. As shown in
Table S15, these changes in the cut-off criteria do not change the signs of our estimates and have
little effect on their magnitudes.

1491 Table S15. Sensitivity Analyses: Comparing inferences when using different cutoffs for 1492 defining species as rare or non-rare based on their relative cover at a site. We compare our 1493 results presented in Figure 5 (column 1 with Cut off 1) to two additional cut offs for classifying a 1494 rare versus non-rare species (in columns 2 and 3); these cutoffs are described in Section S8ciii. 1495 Again, we estimate Equation (S1) with species richness disaggregated into the numbers of non-1496 rare native species, rare native species, non-rare non-native species, and rare non-native species 1497 with groups defined based on each of the three cut offs. All estimated effects of each category of 1498 species richness (SR) are on a log-inverse-hyperbolic-sine scale. Given the inverse hyperbolic 1499 sine transformation of the species richness variables, the estimated effects cannot be interpreted as elasticities without further manipulation, but their signs and relative magnitudes can be 1500 compared to each other. Clustered-robust standard errors are used and clustered at the plot level. 1501 1502 To see 95% confidence intervals as well, see the project page (output/TableS15 R CI.tex). 1502

	Cut off 1 [Main Text]	Cut off 2	Cut off 3
Non-rare, Native SR	0.0488	0.0505	0.0097
	(0.0721)	(0.0848)	(0.0717)
Non-Rare, Non-Native SR	-0.1721***	-0.2143***	-0.1746***
	(0.0649)	(0.0736)	(0.0664)
Rare, Non-Native SR	0.0397	-0.0178	-0.0254
	(0.0711)	(0.0658)	(0.0691)
Rare, Native SR	-0.1473 ***	-0.1050**	-0.0889*
	(0.0459)	(0.0445)	(0.0483)
Num. obs.	1,175	1,175	1,175
Num. plots	146	146	146
R^2 (full model)	0.79	0.78	0.78

1520 Signif. Codes: ***: 0.01, **: 0.05, *:0.1; Robust Standard errors in parentheses (clustered at plot level). 1521

1522

1523 S8d. Variation in each species group

Our analysis exploits observed temporal variation in the data, namely year-to-year changes in each species richness group in a plot. To provide more contextual detail on this variation, we show how changes in the richness of each group of species varies over time by site in Figure S12. Figure S13 breaks this variation down even further by site. These figures imply that the 1528 variation in overall species richness in Figure S2 is driven by changes in rare, native species and



1529 non-rare, non-native species.

1531 **Figure S12. Year-to-year changes in the counts of each species group per plot.** (A) non-

- native, rare species richness; (B) native, non-rare species richness; (C) native, rare species
 richness; and (D) non-native rare species richness.
- 1534

1530





1544 Supplementary Discussion

1545 **S9.** Nine Frequently Asked Questions (FAQ) about Dee et al.

1546 An FAQ is not usually added to an SI, but we found this FAQ to be an effective way to 1547 answer a set of common questions that many readers have and to further clarify how our study 1548 differs from prior studies in ecology.

FAQ#1. Does Dee et al. overturn the popular wisdom from over thirty years of research on the effect of species richness on productivity?

1551 No, but we hope our study gives ecologists a different perspective on the popular wisdom, a

1552 *new approach for conducting empirical ecological research, and fruitful avenues to pursue in*

- 1553 theory and experiments.
- 1554 Regarding the new perspective on the popular wisdom, we believe our study:
- 1555 (i) Estimates the average effect of changes in species richness on productivity when
 1556 richness changes as it does in natural ecosystems, as opposed to when it changes via
 1557 manipulations in experimental systems (see FAQs #2-#4)
- (ii) Accounts for the biological complexity of the ecosystem more comprehensively than
 prior observational designs (see FAQs #6-#8).

1560 Our results are, in fact, broadly consistent with experimental findings of biodiversity 1561 manipulations ('BEF experiments' hereafter). We estimate that an increase in the richness of 1562 species typically planted in BEF experiments - i.e., native, relatively common ('non-rare') 1563 species – has a positive effect on aboveground productivity. However, the Nutrient Network 1564 plots contain many more species than are, and can be feasibly, manipulated in experiments; 1565 namely many more rare and non-native species. More rare species and non-native species are 1566 associated with higher species richness (Fig. S10), and most species in these ecosystems are rare 1567 (Fig. S11) as in most ecosystems (71). These are the species that are changing the most from 1568 year to year (Fig. S12 & S13), and they have a different estimated average effect on productivity 1569 than do the native non-rare species (Figure 5, main text).

1570 Regarding the new approach for ecological research, we believe our approach:

1571 (iii) Makes our causal <u>aspirations</u> transparent ("*Exactly what ecological relationship are*1572 *we trying to estimate?*"). See FAQ #2, #3, and #6.

- 1573(iv)Makes our causal <u>assumptions</u> transparent and motivates them through a combination1574of field knowledge and ecological theory ("What exactly are we assuming when we1575give a causal interpretation to an estimated correlation between richness and1576productivity?"). See FAQ#4, #5, #7, and #8.
- 1577(v)Assesses how changes in our causal assumptions affect our inferences ("How would1578our interpretations change if we use alternative assumptions that may be plausible or1579equally as valid as the assumptions we originally made?"). See FAQ #7 and Figure 31580in manuscript.

Our goal was to <u>build on prior work</u> and advance our collective understanding of the role of biodiversity in ecological functions – not to claim a 'final answer.' We make our assumptions about relationships between the data and our inferences as transparent as possible and explore the implications of these assumptions. This transparency, we hope, makes it easier for ecologists to continue to build on this work, by probing and relaxing these assumptions, and assessing their robustness to new data and methods.

FAQ#2. Why is the estimated average effect in Dee et al. negative while most experimental evidence implies the effect is positive? Is it because the research questions differ; in other words, because the causal effect that experiments aim to estimate and the causal effect that Dee et al. aim to estimate differ?

1591 Yes, differences in research questions could be one reason why the signs of the estimated 1592 effects of species richness on productivity are different. Species richness could change in many 1593 ways, and these different ways are unlikely to have the same effects on ecosystem function. 1594 The typical BEF experiment aims to estimate the expected effect from a change in richness 1595 that arises from a random draw from the pool of species that could grow at the study location and 1596 from the potential values of evenness (i.e., the effect of changing richness independent of 1597 composition attributes). That causal effect is helpful for developing theory because it isolates the 1598 effect of species richness separate from changes in species identities and evenness that may 1599 normally accompany a change in richness.

However, for understanding the ecology of natural systems and processes, particularly when
thinking about the conservation implications of anthropogenic change, we and others (e.g., (77))

1602	argue that the most relevant causal effect of richness on
1603	productivity is the expected effect of a change in species richness
1604	that mimics how richness changes in naturally occurring systems -
1605	more precisely, a change in richness that arises from a random draw
1606	from a weighted conditional distribution of species richness
1607	compositions at different richness values. For example, if dominant
1608	species tend to comprise the majority of species additions or
1609	subtractions when diversity is low, but rare species tend to
1610	comprise the majority of species additions or subtractions when
1611	diversity is high, changes in richness in the study should reflect
1612	those probability distributions (i.e., as a plot or site gets more
1613	diverse, the marginal/incremental species should be more likely to
1614	be rare). Our study aims to estimate this causal effect.
1615	So, in sum, "Yes," the different research aims of the

- 1616 experimental and Dee et al. designs could be a reason for the
- 1617 divergence in results (see also **Box 1**).

After reading this FAQ answer, are you thinking that "composition" is a confounder in Dee et al.? If yes, see **Figure 4** in main text.

1618

- 1619 **FAQ#3. Why is the estimated average effect in Dee et al.**
- 1620 negative while most experimental evidence implies the
- 1621 effect is positive? Is it because the effects of richness on
- 1622 productivity is conditional on species identity (a
- 1623 heterogeneous treatment) and the experiments do not
- 1624 plant the same set of species that are found in natural
- 1625 ecosystems?
- 1626 Yes, the estimated effects of species richness on productivity
- 1627 may differ because the set of species growing in Dee et al. and in
- 1628 *experimental studies differ.*
- 1629 To estimate the average effect of richness on productivity independent of other attributes of
- 1630 diversity, experiments would have to randomize *all* other attributes of diversity (e.g., identities,

Box 1: Spatial and temporal dimensions of diversity's effects on productivity

The causal effects described in FAQ#2 are not sufficiently precise because they lack spatial and temporal dimensions. Over what spatial scale and time horizon are we evaluating the effect of a change in richness? For example, are we referring to the difference in productivity in a patch over one year when the species richness changes from X to Y in a manner that mimics naturally occurring processes? Or are we referring to the difference in productivity in a patch over one hundred years when the species richness changes from X to Y in a manner that mimics naturally occurring processes? We return to this issue in FAQ #5.

1631 densities, relative abundances, traits/functional characteristics). Of course, to do so would be 1632 prohibitively expensive or logistically infeasible. Instead, to our knowledge, most biodiversity-1633 ecosystem functioning (BEF) experiments hold *planted* evenness *constant* at 1, with some 1634 exceptions (for instance, an experiment by Wilsey & Polley (2004) randomize richness across two values of evenness). Some experiments, for a subset of treatments, consider crosses between 1635 1636 functional and species diversity (78). But it's logistically impossible to do all combinations. 1637 Moreover, BEF experiments plant some combinations of species from the set of all possible 1638 species identities and compositions that grow at a site, but not all. For example, in some places, 1639 sourcing seeds for rare species is prohibitively expensive. To our knowledge, BEF experiments 1640 do not plant all combinations of species identities, particularly at higher richness values.

1641 In other words, the species planted in the experimental designs do not comprise the full set of 1642 species that may naturally grow at a study site. Recall that in Dee et al., the negative estimated 1643 effect of richness on productivity appears to come from rare species and non-native species. For 1644 native, dominant species, the estimated effect of richness on productivity is positive. This latter 1645 result from Dee et al result is consistent with BEF experiments, because these native, dominant 1646 species typically comprise most species planted in experiments. Few experiments plant truly rare 1647 species or non-native species, or in the proportions found in natural systems. We further test this 1648 conjecture by analyzing two long-term biodiversity experiments with our study design and find 1649 evidence consistent with the conjecture (see Section S10 for more information).

1650 In the few experiments that have planted rare species, the rare species are a small fraction of 1651 the total species that have been planted in nearly all cases (*but see* the Jena experiment (79)). In 1652 the few experiments that have planted non-native species, the non-native species are also a small 1653 fraction of the total species that have been planted. For example, BioCon includes two non-1654 native species out of the 16 planted, and these two were naturalized to the site (P. Reich, pers 1655 comm). An exception is Wilsey et al. (2009), which paired species according to their native vs 1656 introduced status at a site in Texas, USA (80), where a maximum of 9 non-native species were 1657 planted, out of 20 total planted species.

Even in BEF experiments that plant rare species, the planted rare species may have been less likely to emerge and persist. This phenomenon is known as 'non-compliance' with treatment assignment.³ With non-compliance, *Z* species are planted in a plot (planted richness), but only *D*

³ Non-compliance is also a common, and frequently discussed, challenge in medical randomized control trials.

1661 < Z species emerge or persist (i.e., the 'realized richness') and subsequently affect productivity. 1662 If the species that have a negative or zero effect on productivity are less likely to grow or 1663 survive, it is possible that the true average effect of species richness on productivity in the 1664 experiments could be zero or negative ("true effect" meaning the effect if all planted species were to survive), while the "as planted" estimated average effect of richness on productivity 1665 1666 reported by the experimenters is positive. Analyzing an example of experimental data, we find 1667 evidence for differential non-compliance – some species (e.g., Anemone cylindrica) that are rare 1668 in natural systems in Minnesota, USA do not frequently emerge in the experiment when planted, 1669 particularly in high diversity plots (Box 2). To determine relative abundance and rarity in natural 1670 communities, we use long-term survey data of grasslands around Minnesota, USA from the

- 1671 Minnesota Department of Natural Resources (MN DNR); see dataset details in (81).
- 1672

1673 1674 Are rare species less likely than other species 1675 to be observed growing after being planted in 1676 BEF experiments? To shed light on this 1677 question, we first analyze relative cover and relative frequency data from MN DNR's 1678 1679 grassland monitoring. These data comprise 1680 over 41,000 observations of species, sampled 1681 over the last 10 years from 701 transects 1682 (each with ~25 plots). Using these data, we 1683 then classify species from BioCon 1684 experiments (Cedar Creek) as naturally rare



1685or non-rare (dominant or subordinate) based on the metrics and cut-offs described in the Dee et al SI1686(section S8). Then we calculate the probability that a species is observed growing at each planted1687diversity level (1,2,4,8,16) by comparing the planted versus realized species data. We can then assess1688if "non-compliance" varies based on the species classifications from the MN DNR data (rare,1689subordinate, dominant). We find that rare species are less likely be observed post-planting. Some1690rare species almost never emerge in low diversity plots and never emerge in the 16-species plots1691(e.g., Anemone cylindrica).

1692 So, in contrast to the causal effect that experiments aim to estimate (described in FAQ #2),

1693 the causal effect that experiments actually estimate is the expected effect of a change in richness

1694 that arises from a random draw from a <u>subset of species</u> that could grow at the site, while

holding evenness to a specific value. These causal effects from different subsets, or versions, of
 richness can be helpful for shedding light on the mechanisms through which biodiversity affects
 ecosystem functions -- and thus for shedding light on the heterogeneity of biodiversity's effects
 on ecosystem functions.⁴ But they may not match, in sign or magnitude, the target causal effect
 that experimentalists aim to estimate.

In sum, "Yes," the results from experimental designs and Dee et al.'s design may differbecause different species are planted or survive in the two designs.

FAQ#4. Why is the estimated average effect in Dee et al. negative while most experimental evidence implies the effect is positive? Is it because of potential statistical biases in Dee et al., like unobserved confounding variables or reverse causality, or potential statistical biases in experimental designs?

1706 Yes, biases in either experimental designs or in Dee et al.'s design could create the 1707 divergence in results. In our study, we highlight a variety of potential sources of hidden bias in 1708 their design: for example, unobserved confounders, including measurement error, and reverse 1709 causality, where the effect of productivity on diversity masks or mimics the effect of diversity on 1710 productivity. We attempt to shed light on whether these sources of potential bias could plausibly 1711 mask a true positive relationship between richness and productivity. We find no evidence for 1712 such masking, but, of course, absence of evidence is not evidence of absence. Note that the estimation method itself (i.e., the statistical model) is not likely to yield a negative effect when 1713 1714 the true effect is positive. Indeed, applying our estimation method to experimental data (with a 1715 time-varying treatment or before-after data) gives the same answer as a simpler method that 1716 simply compares mean differences among treated and control groups (section S10). 1717 Experimental designs with randomized treatments can also have hidden biases (reviewed in 1718 (82)). For example, a design that uses comparisons of productivity across different levels of 1719 planted richness to estimate richness's effect on productivity could be biased if randomized 1720 planting, followed by weeding of non-planted species, affects productivity through channels

⁴ However, a broader suite of mechanisms could be operative in natural systems or the relative importance of the mechanisms could differ (92).

other than richness (e.g., a channel like soil disturbance). Wilsey & Polley (2006) noted that
"manipulative experiments have the disadvantage of disturbing soil during plot establishment
(84)."

Although hidden biases in one or both types of designs (ours and experimental) could be one reason for the divergence in results, we believe the two reasons outlined in FAQ #2 and FAQ #3 are more likely: different research questions and different sets of species.

1727 FAQ#5. Models and causal inference: could we use structural equation modeling?

The aim of causal inference is to move from a statistical model to a causal model – or more precisely, to move from a statistical interpretation of a model to a causal interpretation. One cannot make that move without making assumptions, some of which are likely to be untestable. In other words, causal inference cannot rely simply on statistical methods (e.g., structural equation modeling and t-tests). Methods must be complemented by assumptions that determine whether a relationship estimated within a model can be interpreted causally (5, 8, 9). Assumptions are required for all casual inference, whether the data come from an experimental

1735 design or an observational design.

These causal assumptions typically matter more than the statistical methods – an insight that is often summarized in variants of the phrase "design matters more than methods." For example, both structural equation modeling (SEM) and our approach use regression models. In fact, the Dee et al. design could be implemented within an SEM framework. On their own, regression models are simply statistical models without any causal content.

1741 The key innovation of our study is the *design* rather than the estimation method. In other 1742 words, our key contribution is the insight that panel data can be exploited to control for a wider 1743 range of confounding variables than prior studies have achieved. We exert this control in two 1744 ways: (1) control for plot-level fixed effects via a deviation-in-means statistical estimation 1745 procedure; and (2) adding site-by-year dummies to our regression estimator. We could have 1746 exerted this control in multi-level or SEM model by (1) controlling for plot-level fixed effects 1747 via a centering transformation of the data; and (2) adding site-by-year dummies to a SEM model. 1748 Like any causal study, a study that uses SEM requires assumptions for causal interpretation of its 1749 estimates. For SEMs that do not use panel data, a large set of assumptions must be satisfied for 1750 each equation in a SEM to be able to estimate, without bias, its target causal effects. For

example, a SEM with a single year of data cannot address potential confounders for each of the
target causal variables unless either: (1) all potential confounders are measured and in the model;
or (2) the SEM has a valid instrument variable (IV) for each causal variable of interest (variable *Z* in Fig. 1). In other fields, scholars view these assumptions as hard to defend and thus there is
some skepticism about such "all-cause" models (with multiple hypothesis tests, one also needs to
maintain the family-wise Type 1 error rate or control the false discovery rate).

1757

FAQs #6 and #7. Readers of earlier versions of our manuscript have asked us to explain how it is possible that Dee et al. estimate a negative average effect of plot-level richness on productivity while Grace et al. (2016, *Nature*) estimate a positive effect, even though both studies use unmanipulated plots from the Nutrient Network.

The Grace et al. study was a seminal study because it took seriously the complexity of the biological system and used a multivariate approach to quantify relationships among different variables. As noted in the main text, we build on their study and others that followed. Although we cannot know with certainty why the results from our study differ, we briefly summarize our intuition in the main text and, in this FAQ, we describe in more detail the most plausible reasons for the differences.

As noted in FAQ#6, Dee et al. builds on the multivariate advance of Grace et al. Both studies are based on strong theory and field experience about the biology of the systems being studied. Although Dee et al.'s design may look, on the surface, to be less complex than Grace et al.'s SEM, the Dee et al. model tries to address the same ecological complexity. Both designs pose biologically informed hypotheses. It is true that the Grace et al. design aims to test many more hypotheses, i.e., it aims to estimate many more causal relationships than Dee et al. try to estimate (at least 15, by our count). Dee et al. chose to narrow the set of research questions to focus on ruling out rival explanations that arise from spatial and temporal biological complexity in grassland ecosystems. Thus, the key distinction between Grace et al and Dee et al is in the designs - not the methods.

1758

FAQ#6. Why is the estimated average effect of plot richness on plot productivity in Dee et al. negative while in the Grace et al. study it is positive? Is it because the causal effect that the studies aim to estimate differs?

No, the target causal effect in Dee et al. is also one of the target causal effects in Grace et al.
Grace et al. try to estimate more causal effects than Dee et al., but both studies aim to estimate
the average effect of a change in plot-level richness on plot-level productivity. It is true that
Grace et al. and Dee et al. measure productivity differently. Grace et al. measure it as total
biomass (sum of live and dead biomass), whereas Dee et al., and many other studies, measure it
as live biomass. But using total biomass, Dee et al. generate the nearly same estimated negative
effect as they do with live biomass (-0.24, CI [-0.37, -0.11]). See code on project page.

1770Aren't the "changes in species richness" in Grace et al., which uses spatial variation1771in richness to estimate richness's effect on productivity, different from the "changes1772in species richness" in Dee et al., which uses temporal variation in richness? More1773specifically, doesn't Grace et al. measure a "long-run effect" of richness on1774productivity, while Dee et al. measure a "short-run effect?"

We believe this question, which we have received from many people, requires some clarification and elaboration about the two study analysis designs. So, we break down our answer in two parts. First, we clarify that *both* designs use spatial variation, but only Dee et al. uses *both spatial and temporal* variation. Second, we clarify when and whether comparing productivity across plots with different values of richness, as the comparison is done in Grace et al., can provide insights into a "long-run effect" of richness on productivity.

Both studies use spatial variation, but each uses this variation differently.

Our analysis in Dee et al. uses **both** *spatial variation across sites and plots* and *temporal variation across both sites and plots*. In Grace et al., the authors use spatial variation *across sites and plots*; temporal variation is not used because there is only one year of data per plot.

1785 The Dee et al. model eliminates the spatial variation that comes from the "between-plots" 1786 comparisons because we believe those comparisons will yield biased inferences about the 1787 relationship between richness and productivity – hidden bias that comes from unobserved 1788 confounding variables (*see our accompanying Rmarkdown tutorial for elaboration and a visual*). 1789 Yet, even if we were to ignore the potential bias from the between-plots comparisons and use

1790 information from both within and between-plot comparisons and within and between-site

1791 comparisons (called *random effects*" estimator in economics, (17)), we would still get a negative

1792 estimated effect: -0.19 (SE=0.07, p= 0.01, 95% CI [-0.33, -0.04]). Leveraging temporal

1793 variation across sites and plots is a key innovation.

1794 So, the negative estimated effect in Dee et al. does not arise because their design does not use 1795 the between-plot, spatial variation in richness.

1796 Can focusing on spatial variation across plots and sites yield insights into the "long1797 run effect" of richness on productivity?

By "long-run effect," we mean the effect on productivity of a more permanent, or persistent, shift in richness that harnesses long-run mechanisms driven by processes like speciation and evolutionary history occurring over long periods of time -- e.g., *does greater diversity explain why certain ecosystems are more productive than others, holding all other factors constant*?

1802 In Dee et al., the estimated effect in the main design is a short-run effect: the effect of a 1803 change of richness on productivity within ecosystems in a year. We rely on the annual within-1804 plot changes in richness because this variation allows us to control for unobserved plot-level and 1805 site-level confounders that are not easily controlled for in the Grace et al. design with its single 1806 year of data. The annual changes in species richness also allows us to control for reverse 1807 causality in two ways that differ from Grace et al.'s approach (see FAQ#7). We believe that 1808 estimates of the short-run effects of changes in richness on productivity are ecologically relevant 1809 because they capture the effects of changes in today's ecosystems. Nevertheless, using short-1810 term variation in richness may not provide insights into the effects of more persistent shifts in 1811 richness over long time periods.

We next explain and explore considerations for estimating a long-run effect. First, we address this question: "If we <u>only</u> use spatial variation in richness across plots, and ignore the within-plot temporal variation, could we infer the long-run effects of richness on productivity?" To do so, we apply a between-plot estimator to our data. We still obtain an estimate that is negative, albeit small and imprecisely estimated because we only have 2 or 3 plots per site (see STATA code on the release of the project page). 1818 The challenge of estimating long-run effects without long-run data is not unique to ecology – 1819 it's a challenge in all empirical science. An important example is the debate over the effects of 1820 climate change, a long-run phenomenon, versus the effects of weather, a short-run phenomenon. 1821 As in the diversity-productivity context, when there are no data at the temporal scale that one 1822 seeks to estimate an effect, one is stuck either (a) drawing inferences about long-run effects 1823 (climate change, persistent changes in biodiversity) by making strong, and hard to justify, 1824 assumptions about the data-generating process, or (b) drawing inferences about short-run effects 1825 (weather, short-run changes in biodiversity) by making more credible assumptions, which come 1826 at the cost of less certainty over the generalizability of the estimated effects to longer time scales. 1827 For these reasons, we opt to estimate the short-run effect.

FAQ#7. Why is the estimated average effect in Dee et al. negative while in the Grace et al. study it is positive? Is it because each study makes different assumptions about what is driving changes in richness and productivity?

1831 Yes, we believe this reason is the main reason for the different results. The key insight from 1832 Dee et al is that imposing different causal assumptions leads to different causal models, and 1833 different models can yield different conclusions. If the underlying assumptions were false, the 1834 estimated correlation between richness and productivity may not reflect a causal relationship. 1835 Here we try to contrast the assumptions being made in the two studies – our main design 1836 (Figure 2) and the Grace et al design. In doing so, we are not criticizing the Grace et al. design, 1837 but rather showing how their assumptions differ and thus can lead to different conclusions. In 1838 this FAQ, we only review the assumptions made in the Dee et al. main design. One strength of 1839 our study is that we also use other designs that require different assumptions for drawing casual 1840 inferences from the data (Figure 3). Through those designs, we probe the robustness of our 1841 results to violations in the assumptions in the main design. The alternative assumptions in these 1842 complementary analyses are described in detail in the SI below (section S5).

1843 Grace et al. assumptions

1844 To estimate the effect of plot-level richness on plot-level productivity using the Grace et al.

- 1845 design, one must assume the following:
- 1846G1. The plot-level "soil suitability" variable, which is a weighted combination of the percent1847silt and percent sand, only affects plot productivity via its effect on plot richness; i.e., soil

1848 suitability is correlated with richness but uncorrelated with the error term in the plot 1849 productivity equation, after conditioning on site productivity. The key insight in Grace et 1850 al. design is that, to estimate the effect of plot richness on plot productivity without bias 1851 (i.e., to control for omitted variables and for reverse causality), one needs a variable that 1852 is correlated with richness, but not correlated with productivity, except through its 1853 correlation with richness; a so-called "instrumental variable" (IV). If, however, soil 1854 suitability affects productivity through channels other than species richness (after 1855 conditioning on other variables in the model), the estimated correlation between richness 1856 and productivity may not reflect a causal relationship.

- 1857a. If assumption G1 were wrong, the SEM in Grace et al. can still control for1858confounders (but not reverse causality) if the only confounders are the three1859observable variables in the plot productivity equation, or if one were willing to1860make a very strong assumption about the covariance between unobserved richness1861and productivity shocks (i.e., confounders).
- G2. There is no non-classical measurement error⁵ in any of the model variables that creates
 bias in the design. Measurement error could come from either field measurements of the
 variables in the model or the imputed missing soil variable values. For example, one
 must assume:
- 1866

- a. No measurement error in the observable control variables in the regression that are also assumed to be correlated with richness (e.g., soil suitability).
- b. No measurement error in productivity that is correlated with richness.
- 1869 G3. Implicit in the design are also assumptions about the nature of heterogeneous causal
- 1870 effects (i.e., the variability across plots in the effect on productivity from changing
- 1871 richness from *X* to *Y* species) and how the moderators of those heterogeneous effects are
- 1872 distributed across plots. Because those assumptions are more technical, we do not discuss
- 1873 them in detail here, but they are described in (85).
- 1874 In summary, Grace et al. recognize that, in a design with only one year of data, control for 1875 unobserved confounders is challenging unless one has a valid IV. A valid IV also controls for

⁵ Classical measurement error in a variable is when the variable is measured with error but that error is independent of the variables value. Non-classical measurement error is when measurement error of the variable is correlated with the variables value.

1876 reverse causality. There is no empirical test that can validate the assumption that there is no

1877 correlation between an IV (soil suitability) and the outcome variable (productivity) except

1878 through the causal variable (richness). One must use theory and field experience to assess its

1879 validity. Other implicit assumptions made in using SEMs can be found in (86, 87).

1880 **Dee et al. assumptions**

1881 With data that varies over space and time, the challenges to inferring causality from 1882 observational data are less formidable than the challenges when using data that only varies across 1883 space. Using panel data, Dee et al. address bias from confounding variables (omitted variable 1884 bias) and bias from reverse causality in multiple and separate ways (see Fig 2 & 3 in 1885 manuscript). The assumptions used in the Main Design are in the manuscript and SI, but we 1886 summarize them here (other implicit assumptions in our main design are reviewed in (17, 88)): 1887 D1. No plot-level confounders that vary over time. We try several approaches to assess the 1888 robustness of results to this assumption and also use alternative designs that do not require

1889this assumption (see Figure 3 and section S7):

a. We test the sensitivity of their inferences to changes in this assumption.

- b. We use an alternative specification that controls for time-varying, plot-levelconfounders that are correlated with plot productivity, lagged one year.
- c. To address reverse causality, we use an instrumental variable (IV) design (a different
 IV from the IV used in Grace et al; see assumption D3). This design also controls for
 time-varying plot-level confounders when the IV is valid.
- D2. No non-classical measurement error that would create bias in the design (89). The only
 "control variables" are site and year, and thus one may reasonably assume that these
 variables are measured without error. Thus, the only measurement error of concern is error
 in productivity measures that are correlated with richness (or correlated with the predictors
- 1900 of richness in the instrumental variable design).
- D3. In the IV design, the species richness of neighboring plots (the experimentally manipulated plots) only affects plot productivity via its effect on plot species richness; i.e., richness of neighbors is correlated with own richness but uncorrelated with the error (disturbance) term in the own plot productivity equation after conditioning on time-invariant plot attributes and time-varying site attributes. To control for reverse causality, another analysis makes a different set of assumptions than the IV design. We call this the "mechanism blocking" in
 - SI-70

- 1907 Figure 3, and it assumes that one important mechanism through which productivity affects1908 richness is shading (a proposed mechanism from Grace et al.)
- 1909 D4. A weaker version of assumption G3 from the Grace et al design (weaker because we do1910 not rely on a large set of plot or site-level covariates in our estimation procedure).

1911 No empirical test can validate these four assumptions. One must use theory and field experience 1912 to assess their validity and one must probe the robustness of the results to potential violations in 1913 the assumptions. To probe the robustness of our results, we use four approaches (Fig. 3). In each 1914 approach, like our Main Design, causal inference requires causal assumptions. But each 1915 approach makes different assumptions. In other words, each approach can detect hidden biases in 1916 our design (i.e., threats to internal validity) under different conditions. Although each analysis 1917 uses different causal assumptions, each comes to the same conclusion. That pattern is the source 1918 of our paper's strength of evidence: the results of each individual analysis could be explained by 1919 a different rival explanation, but it is harder to come up with a coherent set of rival explanations 1920 that could explain them all (including the results in Fig. 5).

1921

Why would these different assumptions lead to estimated effects of opposite sign?

Here, we elaborate on the logic that underlies the progression of analyses in the Results section of the main text. This progression also highlights how our study builds on prior studies. Let's start by assuming that the data-generating process in the Nutrient Network grassland ecosystems is such that the true average effect of an increase in species richness on live biomass is negative (perhaps small).

1927 Could the three designs generate different conclusions from the same data set? Before using
1928 Nutrient Network data to illustrate how the designs can indeed generate different conclusions, we
1929 explain the intuition via four arguments:

In an observational design using Nutrient Network data, there are plot and site attributes
 that affect richness and productivity in the same directions and that affect richness and
 productivity in opposite directions (e.g., nitrogen positively affects productivity and
 negatively affects richness). If these positive and negative confounding forces are roughly
 similar in magnitude, the simple bivariate correlation between richness and productivity

would be close to zero and thus would be difficult to detect statistically. The estimatedcorrelation may look weakly positive or negative, depending on the sample.

- The plot and site-level confounders that are typically measured in ecological field work
 tend to move richness and productivity in opposite directions, on average (e.g., nitrogen
 negatively affects richness and positively affects productivity). Controlling for <u>only</u> these
 variables thus yields a positive estimated correlation between richness and productivity.
- 1941 3. The plot and site-level confounders that are typically <u>not</u> collected in ecological field work
 1942 tend to move richness and productivity in the same directions, on average (e.g., growing
 1943 season precipitation). If one could observe these variables and only controlled for these
 1944 variables, the estimated correlation between richness and productivity would be negative
 1945 and much larger than the true negative effect.
- 4. Controlling for both the observed and unobserved sources of bias leads to a negative
 estimated effect of species richness on productivity, but one that is smaller than the
 estimate from (3). Arguments (1)-(3) reflect what is sometimes called "point-by-point
 bias," which arises when one controls for some sources of statistical bias but not others.
 "Point-by-point bias" results in the estimated effect moving further from the true effect in
 comparison to when there are no controls for statistical bias.

1952 Now, using the Nutrient Network data, we present three analyses that yield patterns consistent1953 with the four arguments:

- 1954A. Bivariate Correlation between Richness and Productivity: If we were to do an Adler et al-1955like analysis (controlling for year, given we have multiple years), we obtain an Adler et al-1956like result: the estimated correlation is positive (b = 0.14) but we cannot reject the null1957hypothesis that the correlation between richness and productivity is zero (p = 0.17). If we1958just used data from a single year, we find a positive, but statistically insignificant, estimated1959effect in some years and a negative, but statistically insignificant, effect in other years.
- B. Multivariate Model that Controls for Observed Confounders: Multivariate models not only
 have higher explanatory power than traditional bivariate analyses, they also can control for
 confounders when the confounding variables are observable. In Grace et al., a multivariate
 model has much higher explanatory power than the bivariate model and yields a larger,
 positive effect of richness on productivity than the bivariate model. Here, we construct a
 multivariate model that yields similar results. We assume that we do not have a valid IV,
1966 but we can try to control for many potential confounders using Nutrient Network data. If 1967 we were to build on the bivariate analysis via a multivariate analysis that controls just for 1968 17 soil variables, we obtain a result consistent with prior multivariate analyses: a positive 1969 estimated effect (b=0.28) for which one can reject the null hypothesis of the correlation 1970 being equal to zero (p=0.04). Another similar result: in contrast to the low R² in the bivariate analysis (0.04), the R^2 of multivariate analysis is much higher (0.39). If we add to 1971 1972 the model more possible confounders like weather, country, habitat type, and prior use (60 1973 variables in total), the estimated effect is larger (b=0.38; p=0.02) and the R^2 is 0.55 (that's 1974 an overall R^2 for variation both within and between plots; the between-plot R^2 is 0.95).

1975 C. Multivariate Model that Controls for a Wider Range of Observed and Unobserved 1976 Confounders (the "Common Design in Ecology" presented in this paper): Drawing on our 1977 understanding of the biological complexity of these grassland systems, we know that there 1978 are *many* unmeasured variables that could affect both biodiversity and productivity and 1979 thus could be confounding the relationship between richness and productivity; e.g., land-1980 use intensity and history (e.g., (90)), grazing intensity, pollinator diversity at the site, 1981 drought or extreme precipitation events, annual growing season start time, etc. Even 1982 adjusting for 60 covariates, as in (B), is unlikely to eliminate all confounding effects. 1983 Because we have data over both space and time, we can control for a wide range of 1984 confounders, whether the confounding variables can be observed or not. For this reason, we 1985 build on the multivariate analysis in (B) by first controlling for site-level confounders via 1986 site-by-year variables in the model. The estimated correlation between richness and 1987 productivity is of a similar magnitude to the estimate in (B), but negative (b=-0.21). When 1988 we further control for plot-level confounders, we get the estimated negative effect in Dee et 1989 al. (b = -0.24).

1990

What's happening? The Nutrient Network sites are likely to experience site-specific
"shocks" that vary each year (e.g., weather shocks, like a particularly dry April, or herbivory
shocks, like higher herbivore pressure than the prior year). We don't know what exactly these
shocks are, but because we observe the same sites over many years, we can control for them. In
our data, we observe that these shocks affect productivity and richness in the same direction, on
average. We also observe that many of the observable variables collected by Nutrient Network
researchers affect productivity and richness in opposite directions, on average. So, as noted by

have roughly the same number of sites: 43 in Dee et al. vs 39 in Grace et al.). Grace et al. has
more plots, whereas Dee et al. has more years. Given that the unmanipulated plots in Nutrient
Network are a random sample of plots from the sites, the two samples have similar expected
values for plot and site attributes ("similar" because the sets of sites are not identical). In other
words, the external validity of the two data sets is roughly the same, with Dee et al. perhaps
having an advantage with more sites and more years. Moreover, neither study suffers from low
statistical power (Box 3).

Dee et al. negative while in the Grace et al. study it is positive? Is it because Grace et al. have a larger data set? 2018 No, both studies have large data sets. Both studies 2019 have roughly the same number of sites: 43 in Dee .). Grace et al. has

- 2015 2016 2017
- 2014 FAQ#8. Why is the estimated average effect in

2013 assumptions may be plausible.

decide. Future research could focus on assessing the ecological conditions under which each set of

Grace et al., when one looks at a single year of data, the

bivariate correlation will not show anything of note. If

one does not have a valid IV, then when one controls for

observable confounders (like soil attributes), one will see

a positive correlation between richness and productivity.

However, when one controls for a much larger range of

sources of positive and negative bias, as done in Dee et

So, "Yes," we believe that some or all the difference

attributable to the studies making different assumptions.

al., you see an overall negative average effect.

in results between Dee et al. and Grace et al. is

Which, if any, of these sets of assumptions better

2010 approximates the truth is something a reader must

- 2011

- 2012

1998

1999

2000

2001

2002

2003

2004

2005

2006

2007

2008

2009

2020

2021

2022

2023

2024

2025

2026

SI-74

Box 3: Statistical power

Dee et al. have two or three plots per 43 sites, on average, with at least five years of data per plot. Grace et al. have larger number of plots from 39 sites from a single year. Grace et al.'s larger number of plots helps increase the statistical power of their between-plot design. But Dee et al. have five or more time periods per plot and that helps to increase the statistical power of their within-plot design. The Dee et al. design is underpowered if, like Grace et al., we were to rely only on the between-plot variation in richness. The Grace et al. design is underpowered if, like Dee et al., they were to try to control for all site differences with site-level dummy variables.

FAQ#9. Won't the answer always be that the effect of species richness onproductivity 'depends'?

2029 Yes and no. Although the sign and strength of the relationship between richness and 2030 productivity will likely vary with context, the average effect of changes in species richness on 2031 productivity in an ecosystem is relevant both for science and for practice. Nevertheless, we 2032 acknowledge that elucidating how the relationship between richness and productivity varies is 2033 important. In our study, we make some advances by exploring two sources of heterogeneity: 2034 compositional variations in the construct "richness" (i.e., multiple versions of richness depending 2035 on composition) and variations in site and plot-level moderators (i.e., attributes that are off the 2036 causal path between richness and productivity, like precipitation or clay content of soil, but 2037 which moderate the mechanism effects between richness and productivity). In other words, we 2038 explore the implications of both heterogeneous treatments and heterogeneous average causal 2039 responses (see (91) for more on these concepts).

Regarding richness as a heterogeneous treatment, we demonstrate that the effect of richness on productivity depends on what type of species are changing at a site. We focus on the heterogeneous effects of changes in rare species and dominant species, as well as changes in native and non-native species. But other species attributes may also matter.

2044 Regarding heterogeneous average causal responses from variations in moderating site and 2045 plot conditions, we present some initial results on those moderators, but more work is needed 2046 (see Section S6c & S6d, Tables S3-S6). Estimating heterogeneous treatment effects is 2047 challenging because the probability of false positives grows dramatically as one explores various 2048 interaction effects. Investigating heterogeneity in a treatment effect can quickly become an 2049 exploratory, data-mining exercise in which we look for variation in the estimated effect of 2050 richness conditional on a wide range of site and plot attributes. Instead, we advocate for a focus 2051 to a test of a hypothesis implied by theory. Thus, more detailed exploration of heterogeneous 2052 treatment effects is beyond the scope of the Dee et al. study.

- 2053
- 2054
- 2055
- 2056
- 2057

2058 **S10. Supplementary References**

- S. L. Morgan, C. Winship, *Counterfactuals and causal inference* (Cambridge University Press, 2015).
- 2061 2. D. Rubin, Causal inference using potential outcomes: Design, Modeling, Decisions. J. Am.
 2062 Stat. Assoc. 100, 322–331 (2005).
- 3. J. D. Angrist, J. Pischke, *Mostly Harmless Econometrics: An Empiricist's Companion*(Princeton University Press, Princeton, NJ, 2009).
- 2065 4. G. W. Imbens, J. M. Wooldridge, Recent Developments in the Econometrics of Program
 2066 Evaluation. *J. Econ. Lit.* 47, 5–86 (2009).
- 2067 5. J. Pearl, Causality: Models, reasoning, and inference, second edition (2011).
- 2068 6. D. B. Rubin, Estimating causal effects of treatments in randomized and nonrandomized
 2069 studies. J. Educ. Psychol. 66 (1974), doi:10.1037/h0037350.
- 2070 7. A. E. Larsen, K. Meng, B. E. Kendall, Causal Analysis in Control-Impact Ecological
 2071 Studies with Observational Data. *Methods Ecol. Evol.*, 2041–210X.13190 (2019).
- 2072 8. D. Rubin, Causal inference using potential outcomes: Design, Modeling, Decisions. J. Am.
 2073 Stat. Assoc. 100, 322–331 (2005).
- 2074 9. J. Pearl, Causal inference in statistics: An overview. *Stat. Surv.* **3**, 96–146 (2009).
- 2075 10. G. T. Smith, On Construct Validity: Issues of Method and Measurement. *Psychol. Assess.*2076 17, 396–408 (2005).
- 2077 11. M. E. Strauss, G. T. Smith, Construct validity: Advances in theory and methodology.
 2078 *Annu. Rev. Clin. Psychol.* 5, 1–25 (2009).
- 2079 12. E. Everts, Identifying a particular family humor style: A sociolinguistic discourse analysis.
 2080 *Humor.* 16, 369–412 (2003).
- 2081 13. P. J. Ferraro, M. M. Hanauer, Advances in Measuring the Environmental and Social
 2082 Impacts of Environmental Programs. *Annu. Rev. Environ. Resour.* 39, 495–517 (2014).
- 2083 14. P. R. Armsworth, K. J. Gaston, N. D. Hanley, R. J. Ruffell, Contrasting approaches to

- statistical regression in ecology and economics: FORUM. J. Appl. Ecol. 46 (2009),
 doi:10.1111/j.1365-2664.2009.01628.x.
- V. Butsic, D. J. Lewis, V. C. Radeloff, M. Baumann, T. Kuemmerle, Quasi-experimental
 methods enable stronger inferences from observational data in ecology. *Basic Appl. Ecol.* **19** (2017), doi:10.1016/j.baae.2017.01.005.
- B. M. Bolker, M. E. Brooks, C. J. Clark, S. W. Geange, J. R. Poulsen, M. H. H. Stevens,
 J.-S. S. White, Generalized linear mixed models: a practical guide for ecology and
 evolution. *Trends Ecol. Evol.* 24, 127–135 (2009).
- 2092 17. J. M. Wooldridge, Econometric Analysis of Cross Section and Panel Data.
 2093 Booksgooglecom (2002), doi:10.1515/humr.2003.021.
- 2094 18. P. J. Ferraro, M. M. Hanauer, Advances in Measuring the Environmental and Social
 2095 Impacts of Environmental Programs. *Annu. Rev. Environ. Resour.* 39, 495–517 (2014).
- 2096 19. M. A. Huston, Hidden treatments in ecological experiments: Re-evaluating the ecosystem
 2097 function of biodiversity. *Oecologia*. 110, 449–460 (1997).
- 2098 20. G. W. Imbens, Instrumental Variables: An Econometrician's Perspective. *Stat. Sci.* 29, 323–358 (2014).
- 21. J. D. Angrist, G. W. Imbens, D. B. Rubin, J. D. Angrist, G. W. Imbens, D. B. Rubin,
 Identification of Causal Effects Using Instrumental Variables Linked references are
 available on JSTOR for this article : Identification of Causal Effects Using Instrumental
 Variables. J. Am. Stat. Assoc. 91, 444–455 (1996).
- 2104 22. B. E. Kendall, A statistical symphony: Instrumental variables reveal causality and control
 2105 measurement error. *Ecol. Stat. Contemp. Theory Appl.*, 149–167 (2015).
- 2106 23. A. J. Macdonald, E. A. Mordecai, Erratum: Amazon deforestation drives malaria
- 2107 transmission, and malaria burden reduces forest clearing (Proceedings of the National
- 2108 Academy of Sciences of the United States of America (2019) 116 (22212-22218) DOI:
- 2109 10.1073/pnas.1905315116). Proc. Natl. Acad. Sci. U. S. A. 117, 20335 (2020).
- 2110 24. A. J. MacDonald, A. E. Larsen, A. J. Plantinga, Missing the people for the trees:

- Identifying coupled natural-human system feedbacks driving the ecology of Lyme
 disease. J. Appl. Ecol. 56, 354–364 (2019).
- 2113 25. H. S. Wauchope, T. Amano, J. Geldmann, A. Johnston, B. I. Simmons, W. J. Sutherland,
 2114 J. P. G. Jones, Evaluating Impact Using Time-Series Data. *Trends Ecol. Evol.* 36 (2021),
 2115 pp. 196–205.
- 2116 26. J. Fieberg, M. Ditmer, Understanding the causes and consequences of animal movement:
 2117 A cautionary note on fitting and interpreting regression models with time-dependent
 2118 covariates. *Methods Ecol. Evol.* (2012), doi:10.1111/j.2041-210X.2012.00239.x.
- 2119 27. D. R. Schoolmaster, J. B. Grace, E. W. Schweiger, B. R. Mitchell, G. R. Guntenspergen,
 2120 A causal examination of the effects of confounding factors on multimetric indices. *Ecol.*2121 *Indic.* (2013), doi:10.1016/j.ecolind.2013.01.015.
- 2122 28. D. R. Schoolmaster, C. R. Zirbel, J. P. Cronin, A graphical causal model for resolving
 2123 species identity effects and biodiversity–ecosystem function correlations. *Ecology* (2020),
 2124 doi:10.1002/ecy.3070.
- 2125 29. B. Shipley, Cause and Correlation in Biology A User's Guide to Path Analysis,
 2126 Structural Equations and Causal Inference (2004).
- 30. J. B. Grace, K. M. Irvine, Scientist's guide to developing explanatory statistical models
 using causal analysis principles. *Ecology*. 101, 1–14 (2020).
- 2129 31. J. Pearl, The Causal Foundations of Structural Equation Modeling. *Handb. Struct. Equ.*2130 *Model.* (2014).
- 2131 32. E. T. Borer, W. S. Harpole, P. B. Adler, E. M. Lind, J. L. Orrock, E. W. Seabloom, M. D.
 2132 Smith, Finding generality in ecology: A model for globally distributed experiments.
 2133 *Methods Ecol. Evol.* 5, 65–73 (2014).
- 2134 33. E. T. Borer, J. B. Grace, W. S. Harpole, A. S. MacDougall, E. W. Seabloom, A decade of
 2135 insights into grassland ecosystem responses to global environmental change. *Nat. Ecol.*2136 *Evol.* 1, 1–7 (2017).
- 2137 34. J. B. Grace, T. M. Anderson, E. W. Seabloom, E. T. Borer, P. B. Adler, W. S. Harpole, Y.

2138		Hautier, H. Hillebrand, E. M. Lind, M. Pärtel, J. D. Bakker, Y. M. Buckley, M. J.
2139		Crawley, E. I. Damschen, K. F. Davies, P. A. Fay, J. Firn, D. S. Gruner, S. M. Prober, M.
2140		D. Smith, Integrative modelling reveals mechanisms linking productivity and plant species
2141		richness. Nature. 529, 390–393 (2016).
2142	35.	P. B. Adler, E. W. Seabloom, E. T. Borer, H. Hillebrand, Y. Hautier, A. Hector, W. S.
2143		Harpole, L. R. O. Halloran, J. B. Grace, T. M. Anderson, J. D. Bakker, L. a Biederman, C.
2144		S. Brown, Y. M. Buckley, L. B. Calabrese, C. Chu, E. E. Cleland, S. L. Collins,
2145		Productivity Is a Poor Predictor of. Science (80). 1750, 1750–1754 (2011).
2146	36.	J. E. Duffy, C. M. Godwin, B. J. Cardinale, Biodiversity effects in the wild are common
2147		and as strong as key drivers of productivity. Nature. 549, 261–264 (2017).
2148	37.	B. J. Cardinale, K. L. Matulich, D. U. Hooper, J. E. Byrnes, E. Duffy, L. Gamfeldt, P.
2149		Balvanera, M. I. O'Connor, A. Gonzalez, The functional role of producer diversity in
2150		ecosystems. Am. J. Bot. 98, 572–92 (2011).
2151	38.	F. van der Plas, Biodiversity and ecosystem functioning in naturally assembled
2152		communities. Biol. Rev. 94, 1220–1245 (2019).
2153	39.	F. I. Isbell, D. Craven, J. Connolly, M. Loreau, B. Schmid, C. Beierkuhnlein, T. M.
2154		Bezemer, C. Bonin, H. Bruelheide, E. de Luca, A. Ebeling, J. N. Griffin, Q. Guo, Y.
2155		Hautier, A. Hector, A. Jentsch, J. Kreyling, V. Lanta, P. Manning, S. T. Meyer, A. S.
2156		Mori, S. Naeem, P. A. Niklaus, H. W. Polley, P. B. Reich, C. Roscher, E. W. Seabloom,
2157		M. D. Smith, M. P. Thakur, D. Tilman, B. F. Tracy, W. H. van der Putten, J. van Ruijven,
2158		A. Weigelt, W. W. Weisser, B. Wilsey, N. Eisenhauer, Biodiversity increases the
2159		resistance of ecosystem productivity to climate extremes. <i>Nature</i> . 526 , 574–577 (2015).
2160	40.	Lauenroth, W. K., H. W. Hunt, D. M. Switft, J. Singh, Estimating aboveground net
2161		primary production in grasslands: a simulation approach. Ecol Model. 33, 297–314
2162		(1986).
2163	41.	J. C. Moore, in Encyclopedia of Biodiversity (Second Edition), S.A. Levin, Ed. (Elsevier,
2164		Second., 2013), pp. 648-656.
2165	42.	H. W. Jari Oksanen, F. Guillaume Blanchet, Michael Friendly, Roeland Kindt, Pierre

2166 2167 2168		Legendre, Daniel McGlinn, Peter R. Minchin, R. B. O'Hara, Gavin L. Simpson, Peter Solymos, M. Henry H. Stevens, Eduard Szoecs, vegan: Community Ecology Package (2019).
2169 4 2170	3.	M. Bertrand, E. Duflo, S. Mullainathan, How much should we trust differences-in- differences estimates? <i>Q. J. Econ.</i> (2004), , doi:10.1162/003355304772839588.
2171 4 2172	4.	A. C. Cameron, D. L. Miller, A Practitioner's Guide to Cluster- Robust Inference. <i>J. Hum. Resouces.</i> 50 , 317–372 (2015).
2173 4. 2174	5.	B. Callaway, A. Goodman-Bacon, P. H. C. Sant'Anna, Difference-in-Differences with a Continuous Treatment (2021) (available at http://arxiv.org/abs/2107.02637).
2175 4	6.	S. Gaure, Lfe: Linear group fixed effects. R J. 5, 104–116 (2013).
2176 4 2177	7.	G. Sugihara, R. May, H. Ye, C. Hsieh, E. Deyle, M. Fogarty, S. Munch, Detecting Causality in Complex Ecosystems. <i>Science (80).</i> 338 (2012).
2178 4	8.	P. W. Holland, Statistics and causal inference. J. Am. Stat. Assoc. 81, 945–960 (1986).
2179 4 2180	9.	D. B. Rubin, Causal Inference Using Potential Outcomes. J. Am. Stat. Assoc. 100, 322–331 (2005).
2181 5 2182	50.	K. Kimmel, L. E. Dee, M. L. Avolio, P. J. Ferraro, Causal assumptions and causal inference in ecological experiments. <i>Trends Ecol. Evol.</i> 36 , 1141–1152 (2021).
2183 5 2184	51.	S. Arif, M. A. MacNeil, Predictive models aren't for causal inference. <i>Ecol. Lett.</i> 25 , 1741–1745 (2022).
2185 5 2186	52.	M. F. Bellemare, C. J. Wichman, Elasticities and the Inverse Hyperbolic Sine Transformation. <i>Oxf. Bull. Econ. Stat.</i> (2020), doi:10.1111/obes.12325.
2187 5 2188	53.	B. J. Wilsey, H. W. Polley, Realistically low species evenness does not alter grassland species-richness-productivity relationships. <i>Ecology</i> . 85 , 2693–2700 (2004).
2189 5- 2190 2191	54.	M. A. Leibold, J. M. Chase, S. K. M. Ernest, Community assembly and the functioning of ecosystems: how metacommunity processes alter ecosystems attributes. <i>Ecology</i> . 98 , 909–919 (2017).

- 2192 55. Y. Wang, M. W. Cadotte, Y. Chen, L. H. Fraser, Y. Zhang, F. Huang, S. Luo, N. Shi, M.
 2193 Loreau, Global evidence of positive biodiversity effects on spatial ecosystem stability in
 2194 natural grasslands. *Nat. Commun.* 10, 1–9 (2019).
- 56. J. B. Grace, T. M. Anderson, E. W. Seabloom, E. T. Borer, P. B. Adler, W. S. Harpole, Y.
 Hautier, H. Hillebrand, E. M. Lind, M. Pärtel, J. D. Bakker, Y. M. Buckley, M. J.
- 2197 Crawley, E. I. Damschen, K. F. Davies, P. A. Fay, J. Firn, D. S. Gruner, S. M. Prober, M.
- D. Smith, Integrative modelling reveals mechanisms linking productivity and plant species
 richness. *Nature*. 529, 390–393 (2016).
- A. J. MacDonald, A. E. Larsen, A. J. Plantinga, Missing the people for the trees:
 Identifying coupled natural–human system feedbacks driving the ecology of Lyme
 disease. *J. Appl. Ecol.* 56, 354–364 (2019).
- S. Creel, M. Creel, Density dependence and climate effects in Rocky Mountain elk: An
 application of regression with instrumental variables for population time series with
 sampling error. *J. Anim. Ecol.* 78, 1291–1297 (2009).
- 59. G. W. Imbens, Instrumental Variables: An Econometrician's Perspective. *Stat. Sci.* 29,
 323–358 (2014).
- P. J. Ferraro, J. N. Sanchirico, M. D. Smith, Causal inference in coupled human and
 natural systems. *Proc. Natl. Acad. Sci.*, 201805563 (2018).
- 2210 61. E. T. Borer, E. W. Seabloom, D. S. Gruner, W. S. Harpole, H. Hillebrand, E. M. Lind, P.
- 2211 B. Adler, J. Alberti, T. M. Anderson, J. D. Bakker, L. Biederman, D. Blumenthal, C. S.
- 2212 Brown, L. A. Brudvig, Y. M. Buckley, M. Cadotte, C. Chu, E. E. Cleland, M. J. Crawley,
- 2213 P. Daleo, E. I. Damschen, K. F. Davies, N. M. Decrappeo, G. Du, J. Firn, Y. Hautier, R.
- 2214 W. Heckman, A. Hector, J. Hillerislambers, O. Iribarne, J. A. Klein, J. M. H. Knops, K. J.
- 2215 La Pierre, A. D. B. Leakey, W. Li, A. S. MacDougall, R. L. McCulley, B. A. Melbourne,
- 2216 C. E. Mitchell, J. L. Moore, B. Mortensen, L. R. O'Halloran, J. L. Orrock, J. Pascual, S.
- 2217 M. Prober, D. A. Pyke, A. C. Risch, M. Schuetz, M. D. Smith, C. J. Stevens, L. L.
- 2218 Sullivan, R. J. Williams, P. D. Wragg, J. P. Wright, L. H. Yang, Herbivores and nutrients
- 2219 control grassland plant diversity via light limitation. *Nature* (2014),
- doi:10.1038/nature13144.

- 2221 62. E. W. Seabloom, P. B. Adler, J. Alberti, L. Biederman, Y. M. Buckley, M. W. Cadotte, S. 2222 L. Collins, L. Dee, P. A. Fay, J. Firn, N. Hagenah, W. S. Harpole, Y. Hautier, A. Hector, 2223 S. E. Hobbie, F. Isbell, J. M. H. Knops, K. J. Komatsu, R. Laungani, A. MacDougall, R. 2224 L. McCulley, J. L. Moore, J. W. Morgan, T. Ohlert, S. M. Prober, A. C. Risch, M. 2225 Schuetz, C. J. Stevens, E. T. Borer, Increasing effects of chronic nutrient enrichment on 2226 plant diversity loss and ecosystem productivity over time. *Ecology* (2021), 2227 doi:10.1002/ecy.3218. 2228 63. J. L. M. Olea, C. Pflueger, A Robust Test for Weak Instruments. J. Bus. Econ. Stat. 2229 (2013), doi:10.1080/00401706.2013.806694. 2230 64. P. J. Ferraro, J. N. Sanchirico, M. D. Smith, Causal inference in coupled human and 2231 natural systems. Proc. Natl. Acad. Sci., 201805563 (2018). 2232 65. N. Beck, J. N. Katz, Modeling dynamics in time-series-cross-section political economy data. Annu. Rev. Polit. Sci. 14, 331-352 (2011). 2233 2234 J. G. Altonji, T. E. Elder, C. R. Taber, Selection on Observed and Unobserved Variables: 66. 2235 Assessing the Effectiveness of Catholic Schools. J. Polit. Econ. 113, 151–184 (2005). 2236 67. E. Oster, Unobservable Selection and Coefficient Stability: Theory and Evidence. J. Bus. 2237 *Econ. Stat.* **0**, 1–18 (2017). 2238 68. A. Bell, M. Fairbrother, K. Jones, Fixed and random effects models: making an informed 2239 choice. Qual. Quant. 53, 1051-1074 (2019).
- A. E. Larsen, K. Meng, B. E. Kendall, Causal Analysis in Control-Impact Ecological
 Studies with Observational Data. *Methods Ecol. Evol.*, 2041–210X.13190 (2019).
- 2242 70. P. Allison, Fixed Effects Regression Models (2012).

2243 71. B. J. Enquist, X. Feng, B. Boyle, B. Maitner, E. A. Newman, P. M. Jørgensen, P. R.

- 2244 Roehrdanz, B. M. Thiers, J. R. Burger, R. T. Corlett, T. L. P. Couvreur, G. Dauby, J. C.
- 2245 Donoghue, W. Foden, J. C. Lovett, P. A. Marquet, C. Merow, G. Midgley, N. Morueta-
- Holme, D. M. Neves, A. T. Oliveira-Filho, N. J. B. Kraft, D. S. Park, R. K. Peet, M. Pillet,
- J. M. Serra-Diaz, B. Sandel, M. Schildhauer, I. Šímová, C. Violle, J. J. Wieringa, S. K.
- 2248 Wiser, L. Hannah, J. C. Svenning, B. J. McGill, The commonness of rarity: Global and

- future distribution of rarity across land plants. *Sci. Adv.* **5**, 1–14 (2019).
- 2250 72. F. W. Preston, With respect to commonness or rarity. *Ecology*. 29, 254–283 (1948).
- 2251 73. D. F. Sax, S. D. Gaines, Species invasions and extinction: The future of native
 biodiversity on islands. *Light Evol.* 2, 85–106 (2009).
- 2253 74. S. S. Parker, W. S. Harpole, E. W. Seabloom, Plant species natural abundances are
 2254 determined by their growth and modification of soil resources in monoculture. *Plant Soil*.
 2255 445, 273–287 (2019).
- M. D. Smith, A. K. Knapp, Dominant species maintain ecosystem function with nonrandom species loss. *Ecol. Lett.* 6, 509–517 (2003).
- M. L. Avolio, E. J. Forrestel, C. C. Chang, K. J. La Pierre, K. T. Burghardt, M. D. Smith,
 Demystifying dominant species. *New Phytol.* 223, 1106–1126 (2019).
- D. S. Srivastava, M. Vellend, Biodiversity-ecosystem function research: Is it relevant to
 conservation? *Annu. Rev. Ecol. Evol. Syst.* 36, 267–294 (2005).
- P. B. Reich, D. Tilman, J. Craine, D. Ellsworth, M. G. Tjoelker, J. Knops, D. Wedin, S.
 Naeem, D. Bahauddin, J. Goth, W. Bengtson, T. D. Lee, Do species and functional groups
 differ in acquisition and use of C, N and water under varying atmospheric CO2 and N
 availability regimes? A field test with 16 grassland species. *New Phytol.* 150, 435–448
 (2001).
- A. Weigelt, E. Marquard, V. M. Temperton, C. Roscher, C. Scherber, P. N. Mwangi, S.
 Felten, N. Buchmann, B. Schmid, E.-D. Schulze, W. W. Weisser, The Jena Experiment:
 six years of data from a grassland biodiversity experiment. *Ecology*. 91, 930–931 (2010).
- 80. B. J. Wilsey, T. B. Teaschner, P. P. Daneshgar, F. I. Isbell, H. W. Polley, Biodiversity
 maintenance mechanisms differ between native and novel exotic-dominated communities. *Ecol. Lett.* 12, 432–442 (2009).
- 81. H. Ratcliffe, M. Ahlering, D. Carlson, S. Vacek, A. Allstadt, L. E. Dee, Invasive species
 do not exploit early growing seasons in burned tallgrass prairies. *Ecol. Appl.* (2022),
 doi:10.1002/eap.2641.

- 82. K. Kimmel, L. E. Dee, M. L. Avolio, P. J. Ferraro, Causal assumptions and causal
 inference in ecological experiments. *Trends Ecol. Evol.* (2021), ,
 doi:10.1016/j.tree.2021.08.008.
- B. J. Wilsey, H. Wayne Polley, Aboveground productivity and root-shoot allocation differ
 between native and introduced grass species. *Oecologia*. 150, 300–309 (2006).
- 84. C. D'Antonio, S. E. Hobbie, in *Species invasions: insights into ecology, evolution and biogeography.*, G. S. Sax DF, Stachowicz JJ, Ed. (Sinauer Associates, Inc., Sunderland, MA, 2005), pp. 65–85.
- M. Lechner, A. Strittmatter, Practical procedures to deal with common support problems
 in matching estimation. *Econom. Rev.* 38, 193–207 (2019).
- 86. K. A. Bollen, J. Pearl, *Handbook of Causal Analysis for Social Research* (Springer
 Netherlands, Dordrecht, 2013; http://link.springer.com/10.1007/978-94-007-6094-3), *Handbooks of Sociology and Social Research*.
- 2289 87. P. W. Holland, ETS Res. Rep. Ser., in press, doi:10.1002/j.2330-8516.1988.tb00270.x.
- 2290 88. P. J. Ferraro, J. J. Miranda, Panel Data Designs and Estimators as Substitutes for
 2291 Randomized Controlled Trials in the Evaluation of Public Programs. *J. Assoc. Environ.*2292 *Resour. Econ.* 4, 281–317 (2017).
- E. Battistin, A. Chesher, Treatment effect estimation with covariate measurement error. *J. Econom.* 178 (2014), doi:10.1016/j.jeconom.2013.10.010.
- 2295 90. E. Allan, P. Manning, F. Alt, J. Binkenstein, S. Blaser, N. Blüthgen, S. Böhm, F. Grassein,
 2296 N. Hölzel, V. H. Klaus, T. Kleinebecker, E. K. Morris, Y. Oelmann, D. Prati, S. C.
- 2297 Renner, M. C. Rillig, M. Schaefer, M. Schloter, B. Schmitt, I. Schöning, M. Schrumpf, E.
- 2298 Solly, E. Sorkau, J. Steckel, I. Steffen-Dewenter, B. Stempfhuber, M. Tschapka, C. N.
- 2299 Weiner, W. W. Weisser, M. Werner, C. Westphal, W. Wilcke, M. Fischer, Land use
- intensification alters ecosystem multifunctionality via loss of biodiversity and changes to
 functional composition. *Ecol. Lett.* 18, 834–843 (2015).
- P. J. Ferraro, A. Agrawal, Synthesizing Evidence in Sustainability Science through
 Harmonized Experiments: Community monitoring in common-pool resources. *Proc. Natl.*

- 2304 Acad. Sci. USA (2021).
- 2305 92. U. Brose, H. Hillebrand, Biodiversity and ecosystem functioning in dynamic landscapes.
 2306 *Philos. Trans. R. Soc. B Biol. Sci.* **371** (2016).

2307