

1 **A workflow for processing global datasets: application to** 2 **intercropping**

3 Rémi Mahmoud¹, Pierre Casadebaig^{1*}, Nadine Hilgert², Noémie Gaudio¹

4 (1) AGIR, Univ. Toulouse, INRAE, Castanet-Tolosan, France

5 (2) MISTEA, Univ. Montpellier, INRAE, Institut Agro, Montpellier, France

6 (*) Corresponding author (pierre.casadebaig@inrae.fr)

7 **ORCID iDs**

8 [Rémi Mahmoud - https://orcid.org/0000-0003-0853-0834](https://orcid.org/0000-0003-0853-0834)

9 [Pierre Casadebaig - https://orcid.org/0000-0001-7225-936X](https://orcid.org/0000-0001-7225-936X)

10 [Noémie Gaudio - https://orcid.org/0000-0002-4528-9851](https://orcid.org/0000-0002-4528-9851)

11 **Abstract**

12 Field experiments are a key source of data and knowledge in agricultural research. An
13 emerging practice is to compile the measurements and results of these experiments (rather
14 than the results of publications, as in meta-analysis) into global datasets. Our aim in the
15 present study was to provide several methodological paths related to the design of global
16 datasets. We considered 37 field experiments as the use case for designing a global dataset
17 and illustrated how tidying and disseminating the data are the first steps towards open
18 science practices. We developed a method to identify complete factorial designs within global
19 datasets using tools from graph theory. We discuss the position of global datasets in the
20 continuum between data and knowledge, compared to other approaches such as meta-analysis.
21 We advocate using global datasets more widely in agricultural research.

22 Introduction

23 Field experiments, whether conducted on farms or at experimental research stations, have
24 traditionally been the primary approach for acquiring knowledge in crop sciences (Maat, 2011).
25 Yet, extrapolating applicable principles from localized experiments remains a challenging
26 task (Makowski et al., 2014). To derive general rules about agroecosystem functioning, meta-
27 analysis, *i.e.*, a “statistical analysis of a large collection of analysis results from individual
28 studies to integrate the findings” (Glass, 1976), is typically employed. Alternatively, global
29 datasets, corresponding to the aggregation of observations from numerous experiments, can
30 serve as another valuable tool for analyzing agronomic data. While the use of meta-analysis
31 to report results is growing in crop science, it is not a mainstream analysis method compared
32 to reports based on a repeated (years) set of one or two field trials. Distinguishing themselves
33 from meta-analyses, global datasets compile raw experimental results on a detailed scale,
34 such as repeated measurements on individuals or multiple state variables on the canopy.
35 In contrast, meta-analysis is typically restricted to published results with a limited set of
36 variables.

37 Although examples of comprehensive agronomic datasets exist (Kattge et al., 2011; Newman
38 and Furbank, 2021), only a few studies have been based on global datasets (Licker et al., 2010;
39 Lobell et al., 2020; Newman and Furbank, 2021) with even less focus on methods for this type
40 of datasets in crop science (Senft et al., 2022). One significant advantage of agronomic global
41 datasets relies on the fact that they include diverse phenotypic observations from varying
42 soils and climates, enabling more reliable generalization of local findings (Tardieu, 2020).
43 These datasets reduce the risk of spurious correlations (Krajewski et al., 2015; Tardieu,
44 2020) and maximize the utility of experimental data yet to be used in scientific publications
45 (Zamir, 2013).

46 However, global datasets come with their own challenges. Assembling these datasets requires
47 extensive data collection, standardization, and homogenization across diverse experiments
48 conducted by different research teams (White and Van Evert, 2008; Makowski et al., 2014).
49 This tedious curation step is an undervalued task, whose duration could be reduced from
50 the adoption of good practices upstream. Recent efforts and international initiatives aimed
51 at opening and standardizing data are emerging, highlighting that data standardization
52 is crucial for improving the interpretation of experimental results and the generalization
53 of knowledge acquisition. It also facilitates statistical meta-analysis and data publication
54 (Krajewski et al., 2015). However, datasets for plant and crop measurements in controlled
55 field trials are still scarce in public databases. The different field experiments gathered often
56 have diverse objectives, leading to unbalanced and incomplete designs. Confounding factors,

57 ~~i.e.~~ i.e. the unintended mixing of two or more effects making them indistinguishable, can also
58 be challenging (Casler, 2015). Consequently, using and analyzing global datasets require a
59 thorough understanding of the dataset, judicious interpretation of the results, identification
60 of balanced data subsets for specific research questions, and acceptance that the effects of
61 some factors may remain indistinguishable. Therefore, the application of statistical learning
62 techniques on global datasets is only feasible after extensive data pre-processing.

63 Despite these challenges, crop science would greatly benefit from the study of global datasets
64 combining multiple experiments (White and Van Evert, 2008; Zamir, 2013; Cruz and
65 Nascimento, 2019). This approach is particularly relevant considering the current agricultural
66 landscape, where crop diversification is crucial for sustainable farming (Duru et al., 2015).
67 This diversification mandates extensive experimentation, requiring robust data-federation
68 efforts. The joint analysis of global datasets makes it possible to understand the context-
69 dependent nature of diverse experiments and enhances comprehension of the interaction
70 between crop diversity and agroecosystem functioning.

71 To achieve this, we recommend adopting practices for designing and analyzing global datasets
72 that align with tidy data (Wickham, 2014; Broman and Woo, 2018) and FAIR principles
73 (Findable, Accessible, Interoperable, and Reusable) (Wilkinson et al., 2016). As a use case,
74 we illustrate the design of a global dataset for intercropping systems, in which at least two
75 crop species are grown in the same field for a significant part of their growth cycle. We
76 describe the main steps involved in designing a global dataset gathering 37 intercropping
77 experiments across Europe. We also describe and apply an original method ~~for identifying~~
78 ~~factorial designs, which is to identify complete factorial design subsets of interest.~~ This
79 methodological development was aimed at helping the potential collaborators to explore
80 and get an overview of the dataset as a function of their factor of interest, a key step in
81 assisting further modeling and analysis steps.

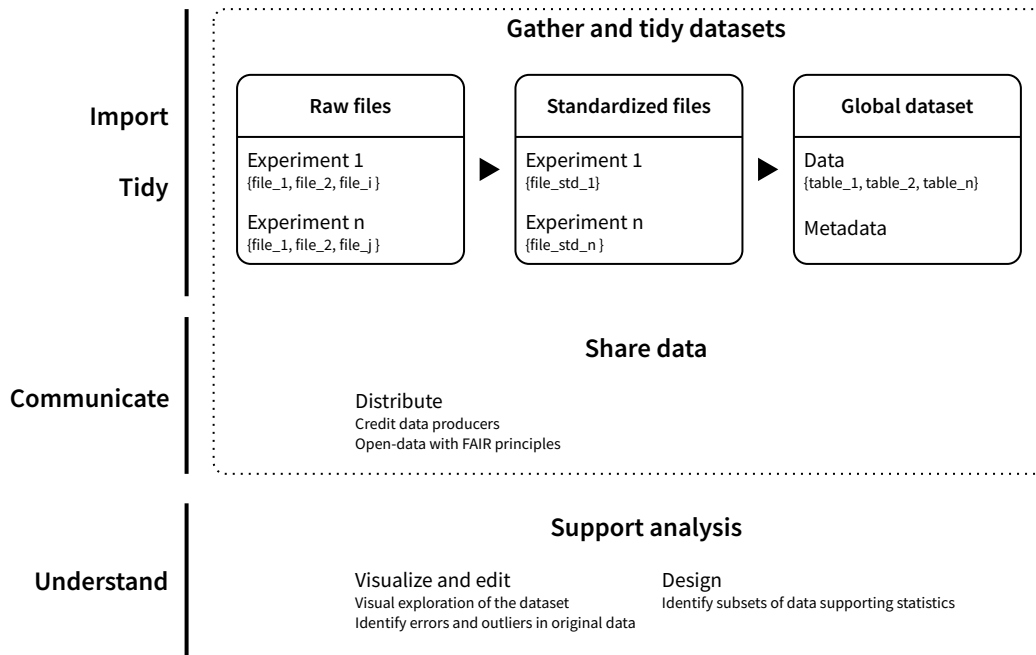
82 Our global aim was to describe our workflow in a realistic manner, hoping to promote these
83 practices and to encourage the scientific community to move towards a more open approach
84 to conducting experimental science in agronomy, making it more reproducible and shared.

85

86 **Design steps of global datasets**

87 This section presents the generic steps involved in designing a global dataset. As the
88 gathering, cleaning, and formatting of the spare source datasets is time-consuming, we

89 followed tidy data specifications (Wickham, 2014) and a global data science workflow as
 90 presented by Wickham and Grolemund (2016) (Figure 1).



91

92 **Figure 1. Main steps for designing global datasets.** The left column corresponds to a
 93 classical data science workflow. We adapted these steps for global dataset design specificities,
 94 to illustrate the importance of data gathering, tidying, and sharing (dotted frame). While some
 95 actions supporting subsequent data analysis are generic (visualization, editing), most depend
 96 on the chosen analysis strategy.

97 1. Gather and tidy source datasets

98 1.1. Conceptual framework

99 Overall, the aim of this gathering and tidying step is to transform a highly heterogeneous
 100 set of tables, scattered in various files according to the logic of each practitioner, into a
 101 structured and documented set of rectangular files.

102 In a first step, the research groups that conducted the experiments whose features are
 103 interesting for a global dataset shall be identified and contacted. While the data processing
 104 step is often known to be very time-consuming in the overall data science workflow (Wickham,
 105 2014), this contact and convincing step is also very long, with potential disappointing
 106 responses (Popkin, 2019).

107 Then, a basic database model for the global dataset has to be developed. This step
108 involves defining the structure of a database, including the number of tables needed and the
109 relationships between them. It also involves describing the metadata, such as the variables
110 measured or collected, their definitions, and units.

111 Using this database model, the raw experimental files are standardized, from various
112 spreadsheet formats into a single and coherent dataset. In crop science, operating by field
113 experiment makes the whole process easier, by focusing standardization efforts on a set
114 of files sharing common properties (illustrated by moving from *raw* to *standardized* files
115 in Figure 1). These standardized files are then combined and documented to make the
116 data “analysis-friendly” (Wilson et al., 2017), which enables detection of errors and data
117 exploration, validation and analysis. A good practice is to work with “tidy” data which is
118 a standard way of mapping the meaning of a dataset to its structure (Wickham, 2014). A
119 dataset is messy or tidy depending on how rows, columns and tables are matched up with
120 observations, variables and types. In tidy data, every column is a variable, every row is an
121 observation, and every cell is a single value. Messy data is any other arrangement of the
122 data (Wickham and Grolemund, 2016; Broman and Woo, 2018).

123 1.2. Case study

124 ~~While there are relatively few incentives to share agronomical (Senft et al., 2022) or~~
125 ~~ecological (Jenkins et al., 2023) datasets, requirements and practices need to evolve. The~~
126 ~~ability to easily disseminates data is thus a key feature in designing a dataset, since it~~
127 ~~determines how other researchers will be able to interact with the data, and potentially~~
128 ~~increase its reuse. Open data should be designed in accordance with the FAIR data~~
129 ~~principles (-).~~

130 ~~When discussing with the involved research groups, one recurrent constraint to open~~
131 ~~their data was the perception that their contribution could not be credited unless sharing~~
132 ~~authorship in research articles. If applied consistently, open data FAIR requirements will~~
133 ~~allow contributors to be specifically acknowledged for their work, through citation of the~~
134 ~~dataset they contributed to (Jenkins et al., 2023).~~

135 ~~Once the data are in a tractable format, visual exploration allows for a comprehensive~~
136 ~~overview of data patterns, aiding in the identification of anomalies such as errors and outliers~~
137 ~~that may not be immediately apparent through numerical analysis alone.~~

138 ~~Later, additional processes are required to render the dataset operational for analytical and~~
139 ~~modeling studies, such as data imputation, dimension reduction, or data normalization.~~

140 ~~Because these steps depend largely on the chosen analytical workflow, they are not directly~~
141 ~~included in the communicated open datasets, but rather tailored by the subsequent~~
142 ~~analytical team (dotted frame in Figure 1).~~

143 ~~Nonetheless, sharing methods can support the future reuse of the dataset. In our case in~~
144 ~~crop ecology, we illustrated this step with the development of an original method aiming at~~
145 ~~identifying subsets in the overall dataset corresponding to complete factorial designs. This~~
146 ~~method is presented in the following section.~~

147 ~~We briefly describe the features of the available field experiments to highlight their richness~~
148 ~~and heterogeneity (see Gaudio et al. (2021) and Mahmoud et al. (2022) for full details and~~
149 ~~experimental protocols; see Gaudio et al. (2023) for the global dataset online).~~

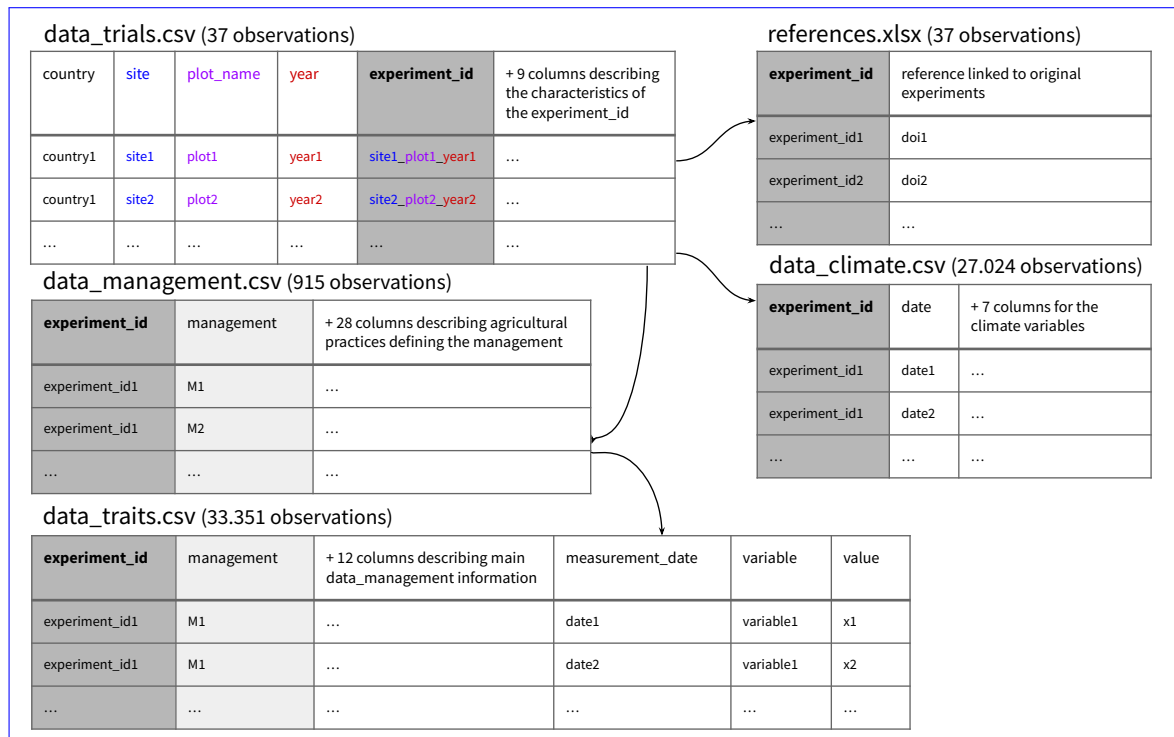
150 Although combining results from a few experiments (usually two years, often sequential) is
151 common in the intercropping literature (and more generally in crop science), no study includes
152 joint analysis of dozens of experiments to infer more generic results about intercropping
153 functioning. To this end, we designed, built and analyzed a global dataset gathering the
154 results of 37 field experiments that involved cereal-legume intercrops and the corresponding
155 sole crops. Globally, the aim of these field experiments was to compare the growth and
156 grain yield ($\text{t}\cdot\text{ha}^{-1}$) of multiple combinations of species grown in intercrop to their sole-crop
157 reference. The field experiments were carried in 5 European countries (France, Denmark,
158 Italy, Germany and England) from 2001 to ~~2017 (Figure 2).~~ 2017. The global dataset
159 included 5 legume species (chickpea, faba bean, lentil, lupin and pea), 3 cereal species (barley,
160 durum wheat and soft wheat) and 8 resulting intercrops, *i.e.* i) barley associated with faba
161 bean, lupin or pea, ii) durum wheat associated with chickpea, faba bean or pea, and iii) soft
162 wheat associated with lentil or pea.

163 ~~**Figure 2.** Location of the 37 intercropping experiments gathered within the global dataset.~~
164

165 To gather the 37 experiments, six research teams were contacted. For each experiment, several
166 ~~excel files~~ spreadsheet files (all in Excel format) were retrieved, ranging from 1 to 10 per
167 experiment. These files differed by the number of ~~spreadsheets~~ sheets they contained, ranging
168 from 1 to 67. We finally collected a total of 86 excel files ~~and (412 spreadsheets~~ sheets). These
169 raw data were highly heterogeneous at all levels, whether concerning the variables (~~e.g.~~ e.g.
170 type, name, unit, measured scale) or the format of the file itself (~~e.g. one spreadsheet~~ e.g.
171 one sheet per date or per variable, different tables on a same ~~spreadsheets~~ sheet, calculations
172 and graphs ~~within raw data files~~ mixed with raw data cells, different languages and encoding
173 format).

174 ~~After the step of gathering , the files were~~ Aiming at improving machine and human
175 readability (Wilson et al., 2017), variable names were chosen to be as explicit as possible. We
176 settled for composite names separated by underscore and containing: as few abbreviations
177 as possible, a reference to the organizational levels (organs: leaf, shoot; individuals: plants;
178 population: crop), and a reference to the variable itself (biomass, number, length). After
179 gathering step, the information of the files was transformed into standardized rectangular
180 data tables, following ~~the tidy format and good practices a~~ *tidy* format (Wickham, 2014;
181 ~~and recommended practices of data organization in spreadsheets~~ (Broman and Woo, 2018),
182 resulting in the creation of one given file per experiment. ~~Each file includes~~ The measured
183 values were not normalized (for *e.g.* spatial field or experimenter effects) as the information
184 on experimental design type and structure was only accessible in very few trials. Each file
185 included 6 ~~spreadsheets, in which the variables and values were placed~~ sheets with one table
186 per sheet, defined as a function of the ~~information category of data~~ they provided (~~e.g.~~ *e.g.*
187 plant functioning, climate, agricultural practices). This step resulted in the creation of 37
188 excel files (vs. 86) and 222 ~~spreadsheets~~ sheets (vs. 412).

189 Finally, all the files were pooled together using R software, ~~with a final table per type of~~
190 ~~variable, i.e. to create one global table per data category, i.e.~~ four tables related respectively
191 to climate, crop measurements, agricultural practices and global information describing the
192 site (Figure 2). Overall, the global dataset contained 308 and 299 statistical individuals
193 (~~i.e. defined as~~ a unique combination of {site * year * management}) in intercrop and sole
194 crop, respectively (Table 1). The number of plant characteristics was much larger (33351
195 observations, among which 12896 were measured in sole crops and 20455 in intercrops), since
196 several variables were measured at the crop scale, sometimes several times during the crop
197 cycle.



198

199

200

201

202

203

204

205

206

207

208

209

210

211

Figure 2. Representation of the relationships between tables identified in the global dataset. Five tables were defined to organize data, all sharing a common identifier (*experiment_id*, which is the concatenation of the *site* *plot* *year* of each experiment). The table *data_trials.csv* provides the main characteristics (*e.g.* latitude/longitude, soil texture) of each site, with one line per experiment (37 observations). The table *data_climate.csv* provides the climate time series during the growing season for each experiment (27,024 observations), retrieved using a gridded API (NASA POWER API, Sparks (2018)). The table *data_management.csv* describes the different agricultural practices used in each experimentation (*e.g.* species grown in sole- or intercrop, genotype, fertilization). The table *data_traits.csv* provides all the plant variables and their value as a function of time (measurement) per management and experiment (33,351 observations). Finally, the table *references.xlsx* provides the initial experimental references linked to each experiment (when existing).

212 **Table 1. Overview of the diversity of the treatments in the global dataset by factors**
 213 **(columns) and experiments (rows).** Within each column, each colored rectangle is a level
 214 of the factor considered. For instance, the two colors for the *Mixing pattern* indicate that the
 215 two species intercropped were sown in alternate rows or within the row; the two colors for the
 216 *Nitrogen (N) fertilization* indicate that the experiment included at least two N-treatments (no
 217 fertilization and N-fertilization, the latter of which may include several amounts of N); regarding
 218 *Species mixture*, the number of colors indicates the number of different species mixtures included
 219 in a given experiment. A rectangle in a given row and column indicates that the corresponding
 220 experiment contains at least one statistical individual with the corresponding factor level.

| Experiment | No. of statistical individuals | No. of variables | Mixing pattern | Species mixture | Nitrogen fertilization |
|-----------------------------|--------------------------------|------------------|----------------|---------------------------|------------------------|
| Taastrup_taastrup_2003 | 6 | 9 | Red | Green | Cyan |
| SanMarco_sanMarco_2004 | 4 | 10 | Red | Green | Cyan |
| SanMarco_sanMarco_2003 | 4 | 10 | Red | Green | Cyan |
| Reading_reading_2003 | 6 | 10 | Red | Green | Cyan |
| Kassel_kassel_2004 | 6 | 10 | Red | Green | Cyan |
| Jynde vad_jyn_2003 | 24 | 11 | Red | Green, Yellow, Orange | Cyan |
| Jynde vad_jyn_2002 | 24 | 12 | Red | Green, Yellow, Orange | Cyan |
| Jynde vad_jyn_2001 | 24 | 12 | Red | Green, Yellow, Orange | Cyan |
| Grignon_inrae_2017 | 19 | 6 | Red | Cyan | Cyan |
| Grignon_inrae_2010 | 16 | 8 | Red | Cyan | Red, Cyan |
| Grignon_inrae_2009 | 15 | 8 | Red | Cyan | Red, Cyan |
| Grignon_inrae_2008 | 27 | 5 | Red | Cyan | Red, Cyan |
| Grignon_inrae_2007 | 30 | 7 | Red | Cyan | Red, Cyan |
| Copenhagen_hbg_2003 | 24 | 10 | Red | Green, Yellow, Orange | Cyan |
| Copenhagen_hbg_2002 | 24 | 11 | Red | Green, Yellow, Orange | Cyan |
| Copenhagen_hbg_2001 | 24 | 12 | Red | Green, Yellow, Orange | Cyan |
| Auz_ZN_2012 | 58 | 24 | Red | Green, Purple, Pink | Cyan |
| Auz_TO_2016 | 86 | 18 | Red | Green | Cyan |
| Auz_TO_2013 | 93 | 24 | Red | Green, Purple, Pink | Red, Cyan |
| Auz_TE_2006 | 13 | 20 | Red | Green, Purple, Pink | Red, Cyan |
| Auz_SGs_2007 | 66 | 23 | Red | Green, Purple, Pink | Red, Cyan |
| Auz_PP_2011 | 20 | 20 | Red | Green, Purple, Pink | Red, Cyan |
| Auz_pk_2011 | 18 | 18 | Red | Green, Purple, Pink | Red, Cyan |
| Auz_marinette_2_2015 | 85 | 13 | Red | Green | Cyan |
| Auz_marinette_1_2015 | 22 | 13 | Red | Green | Cyan |
| Auz_cochard_2010 | 60 | 21 | Red | Green, Blue, Purple, Pink | Red, Cyan |
| Angers_thorigne_2009 | 11 | 12 | Red | Cyan | Red, Cyan |
| Angers_thorigne_2008 | 15 | 14 | Red | Cyan | Red, Cyan |
| Angers_thorigne_2007 | 11 | 12 | Red | Cyan | Red, Cyan |
| Angers_thorigne_2006 | 6 | 8 | Red | Cyan | Red, Cyan |
| Angers_thorigne_2004 | 6 | 10 | Red | Green | Cyan |
| Angers_thorigne_2003 | 6 | 10 | Red | Green | Cyan |
| Angers_jailliere_2008 | 22 | 16 | Red | Cyan | Red, Cyan |
| Angers_jailliere_2007 | 14 | 16 | Red | Cyan | Red, Cyan |
| Angers_fnams_2003 | 12 | 10 | Red | Green | Red, Cyan |
| Angers_fnams_2002 | 4 | 8 | Red | Green | Cyan |
| Angers_brainsurlaution_2011 | 10 | 5 | Red | Cyan | Cyan |

221

222 2. Share organized data

223 While there are relatively few incentives to share agronomical (Senft et al., 2022) or
224 ecological (Jenkins et al., 2023) datasets, requirements and practices need to evolve
225 (Krajewski et al., 2015). The ability to easily disseminates data is thus a key feature in
226 designing a dataset, since it determines how other researchers will be able to interact with
227 the data, and potentially increase its reuse. Open data should be designed in accordance
228 with the FAIR data principles (<https://force11.org/info/the-fair-data-principles/>).

229 When discussing with the involved research groups, one recurrent constraint to open
230 their data was the perception that their contribution could not be credited unless sharing
231 authorship in research articles. If applied consistently, open-data FAIR requirements will
232 allow contributors to be specifically acknowledged for their work, through citation of the
233 dataset they contributed to (Jenkins et al., 2023).

234 This global dataset, as well as the metadata associated, are available on a data repository in
235 a FAIR way (Gaudio et al., 2023). Out of the 37 experiments gathered, 11 have never been
236 valued before.

237 Additional details on experimental designs and management practices are reported in the
238 reference publications for 26 of the 37 experiments (Knudsen et al., 2004; Corre-Hellou et
239 al., 2006; Hauggaard-Nielsen et al., 2008; Hauggaard-Nielsen et al., 2009a; b; Launay et al.,
240 2009; Bedoussac and Justes, 2010a; b; Naudin et al., 2010, 2014; Barillot et al., 2014; Pelzer
241 et al., 2016; Tang et al., 2016; Viguier et al., 2018; Kammoun et al., 2021).

242 3. Support new analysis

243 3.1. Conceptual framework

244 ~~Table 1. Diversity of the treatments in the global dataset by factor (columns) and~~
245 ~~experiment (rows). Within each column, each colored rectangle is a level of the factor~~
246 ~~considered. A rectangle in a given row and column indicates that the corresponding~~
247 ~~experiment contains at least one statistical individual with the corresponding factor level.~~
248 Once the data are in a tractable format, visual exploration allows for a comprehensive
249 overview of data patterns, aiding in the identification of anomalies such as errors and
250 outliers that may not be immediately apparent through numerical analysis alone. Later,
251 additional processes are required to render the dataset operational for analytical and
252 modeling studies, such as data imputation, dimension reduction, or data normalization.
253 Because these steps depend largely on the chosen analytical workflow, they are not directly

254 included in the communicated open datasets, but rather tailored by the subsequent
 255 analytical team (Figure 1). Nonetheless, sharing methods can support the future reuse of
 256 the dataset. In our case in crop ecology, we illustrated this step with the development of
 257 an original method aiming at identifying subsets in the overall dataset corresponding to
 258 complete factorial designs.

259 3.2. Case study

260 Method

261 The brief description of the global dataset revealed the diversity of agronomic situations
 262 considered (Table 1). While the experimental designs ~~had share~~ many similarities (~~e.g.-e.g.~~
 263 species cultivated, agricultural ~~management practices~~), the resulting overall design ~~did not~~
 264 ~~allow an immediate statistical analysis of the global dataset is unbalanced~~. We thus developed
 265 a method to *a posteriori* identify subsets in the global dataset corresponding to complete
 266 factorial designs. This approach can quickly assess whether the dataset is suited to answer a
 267 set of scientific questions, as long as the factors of interest are sufficiently represented in the
 268 global dataset. The role of this method was not to identify potential confounding factors,
 269 which is left for the interpretation of the results of further statistical analysis

270 To identify the largest data subsets associated with complete factorial designs in the global
 271 dataset, we used tools from graph theory (Phillips et al., 2019). In graph theory, a graph G
 272 is a pair $G = (V, E)$ where V is a set of vertices, and E is a set of edges that connect some
 273 of the vertices (Table 2).

274 **Table 22. Definitions in graph theory used in the present study.** ~~—Definitions in~~
 275 ~~graph theory used in the present study (Phillips et al., 2019)~~

| Term | Definition |
|--|--|
| <i>subgraph</i> $\tilde{G} = (\tilde{V}, \tilde{E})$ of a graph $G = (V, E)$ | A graph whose vertex set (\tilde{V}) is included in the vertex set of G (i.e.-i.e. $\tilde{V} \subseteq V$) and whose edge set (\tilde{E}) is included in the edge set of G (i.e. $\tilde{E} \subseteq E$) |
| <i>complete graph</i> | A graph whose vertices are all connected |
| <i>clique</i> of a graph G | A complete subgraph of G |
| <i>maximal clique</i> of a graph G | A clique that cannot be extended by including one more adjacent vertex |
| <i>k-partite graph</i> | A graph that can be partitioned into k nonempty non-empty , vertex-disjoint, edgeless subgraphs |
| <i>k-partite clique</i> or <i>k-clique</i> | A set of vertices that induces a complete k -partite subgraph |

| Term | Definition |
|--|---|
| <i>maximal k-partite clique</i> | A k -clique that cannot be extended by including one more adjacent vertex |

276 Given a set of categorical variables X_1, \dots, X_k , each having values in a discrete set (~~i.e.~~i.e.
277 $\forall i = 1, \dots, k X_i \in \mathcal{A}_i := \{x_{i,1}, \dots, x_{i,j_i}\}$, ($j_i \in \mathbb{N}^*$ denoting the number of levels of variable
278 X_i)), a k -partite graph can be derived by setting $V = \bigcup_{i=1}^k \mathcal{A}_i$ (~~i.e.~~i.e. each level of each
279 factor is a vertex) and $E = \{(x, y) \mid \text{levels } x \text{ and } y \text{ observed together}\}$.

280 A factorial design is complete if, and only if, all possible combinations of the factor levels
281 are present. For a graph $G = (V, E)$, this is equivalent to identifying a subgraph with an
282 edge between each pair of vertices from independent sets (~~i.e.~~i.e. a k -clique). Thus, the
283 challenge of identifying the largest complete factorial designs within a global dataset can be
284 reduced to counting the number of maximal k -cliques in the graph.

285 Phillips et al. (2019) developed the Maximum Multipartite Clique Enumeration (MMCE)
286 algorithm to count the number of maximal multipartite cliques within a k -partite graph.
287 MMCE starts from the observation that if G is k -partite, and if another graph G' is built
288 from G by adding all intrapartite edges, then C is a maximal k -partite clique in G if C is a
289 maximal clique in G' with at least one vertex in each partite set. Thus, the initial question is
290 a matter of a modified problem of maximal clique enumeration, which is a NP -hard problem
291 (Lawler et al., 1980). To address this issue, the MMCE algorithm uses a graph inflation
292 approach, by adding all possible intrapartite edges to G . It then identifies maximal cliques
293 in the inflated graph using a procedure of Bron and Kerbosch (1973) and checks whether
294 the cliques identified cover all of the partite sets. We coded MMCE in the R programming
295 language (<https://github.com/RemiMahmoud/kclique>). Although the problem of identifying
296 maximal k -partite cliques with the maximum number of vertices has also been shown to be
297 NP -hard for any $k \geq 3$ (Phillips et al., 2019), the relatively few vertices ($|V| < 300$) in the
298 global dataset allowed solutions to be found quickly.

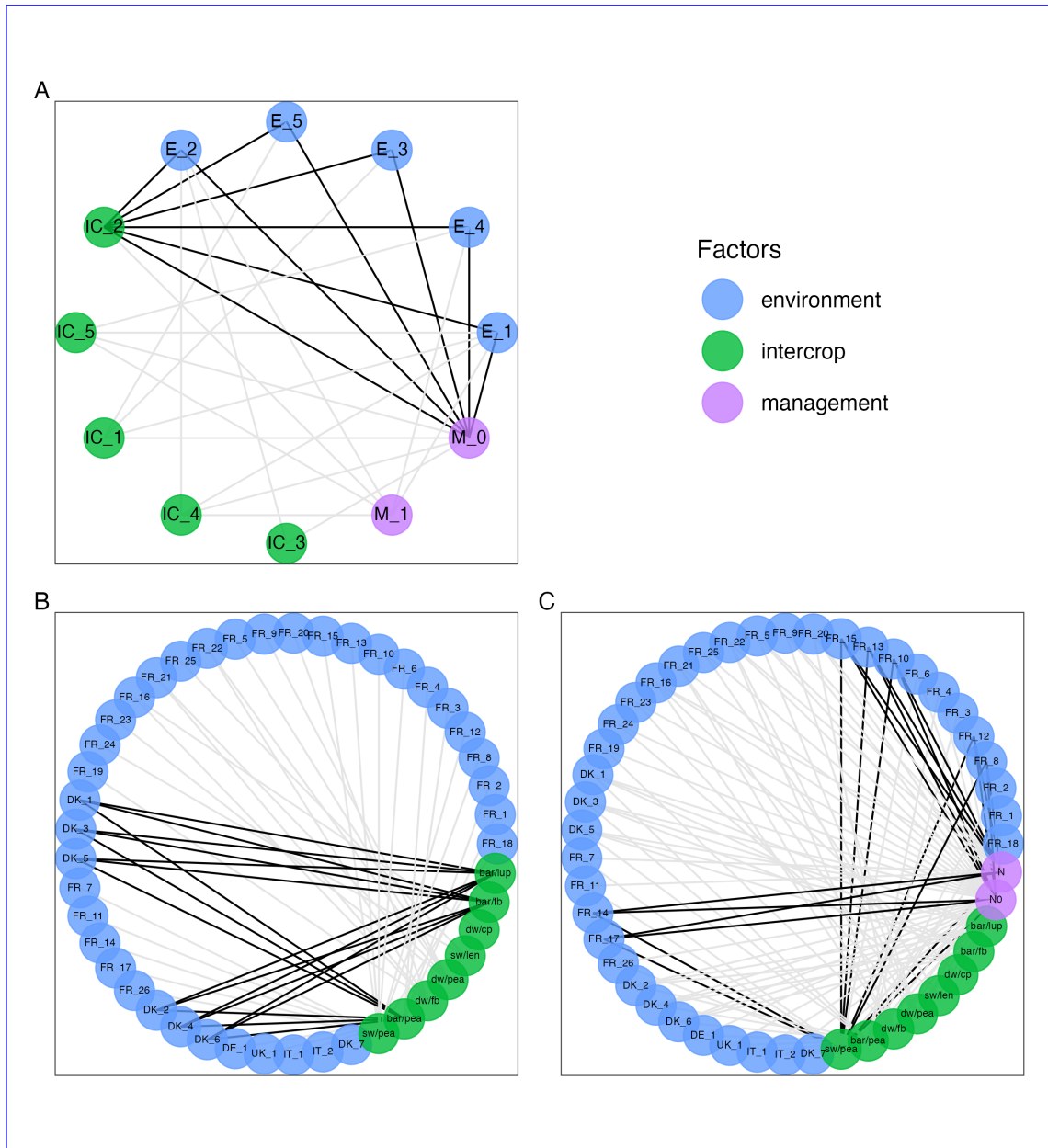
299 Application

300 Here, we illustrate this method with ~~a fictive global dataset made up from two datasets~~ : (1)
301 ~~a theoretical one, where we generated~~ an unbalanced design of five environments (~~site*year~~),
302 ~~five crops~~, ~~five intercrops~~, and two management levels (Figure ~~3~~). ~~3A~~; and (2) a practical
303 ~~one, corresponding to the global dataset presented in this study~~ (Figure 3B and 3C).

304 When applied on ~~this unbalanced design~~ ~~the theoretical unbalanced design~~ (Figure 3A),
305 this method identified ~~11~~ 8 maximal 3-partite cliques, ~~with four examples illustrated in~~
306 ~~Figure 3~~. ~~While~~ each of these ~~examples maximized the representativeness of a factor of~~
307 ~~interest~~ (~~crop~~, ~~environment~~, ~~designs having different number of modalities in considered~~
308 ~~factors~~ (~~environment~~, ~~intercrops~~ or management), ~~There is only one design maximizing~~

309 the number of environments, and no factorial design was found with two levels per factor in
310 ~~this fictive dataset.~~

311 ~~We also applied this method to address a specific issue (Mahmoud et al., 2022), in which we~~
312 ~~analyzed how nitrogen (N) fertilization influenced plant-plant interactions within intercrops.~~
313 ~~To this end, we looked for experiments that included both N-fertilized and unfertilized~~
314 ~~treatments by looking for a maximal 2-clique in a graph composed of two sets of vertices:~~
315 ~~i) field experiments and ii) N fertilization (i.e. unfertilized and N fertilized levels). The~~
316 ~~targeted maximal 2-clique needed to contain the two levels of the sets of N-fertilization~~
317 ~~vertices.~~



318

319 **Figure 3. ~~Four-Three~~ maximal ~~3-cliques~~ ~~k-cliques~~ that represent distinct complete**
 320 **factorial designs within ~~an~~-theoretical (A) and experimental (B-C) unbalanced**
 321 **~~design-with-three factors~~ designs.** Black edges represent the edges of the ~~3-cliques~~ ~~cliques~~
 322 and gray edges represent the factor combinations appearing in the initial design. ~~Despite~~
 323 ~~the potential richness of the global dataset, there was no case where two levels of each factor~~
 324 ~~were combined in a factorial design: network~~ In the case A ~~focused on crops, network,~~ we
 325 ~~generated a random unbalanced design for three factors and illustrated the 3-clique maximizing~~
 326 ~~the number of environments.~~ The experimental design in the cases B ~~on environments, network~~
 327 ~~C on management, and network D on crop and management together.~~ C corresponds to the
 328 ~~aggregation of the 37 experimentations (blue nodes).~~ In case B, we searched for any intercrop

329 observed at least in two environments. In case *C*, there was an additional constraint on two
330 levels of nitrogen (N) fertilization. Countries were abbreviated with their ISO 3166 codes;
331 species were abbreviated as barley (*bar*), chickpea (*cp*), durum wheat (*dw*), faba bean (*fb*),
332 lentil (*len*), lupin (*lup*), soft wheat (*sw*); nitrogen fertilization was abbreviated as *N0* for no
333 fertilization, and *N* for fertilization.

334 We considered two examples for the application on the agronomic global dataset. In the
 335 first one, we searched for any number of intercroops observed at least in two environments.
 336 Two designs were identified: the one with the most environmental modalities is illustrated
 337 in Figure 3B; the alternative design was, crossing {environments} x {intercroops}, {FR 22,
 338 FR 21} x {dw/pea, dw/fb}. The second example was the same request with an additional
 339 constraint on two levels of nitrogen (N) fertilization. In this case, three designs were
 340 identified, the largest one being illustrated in Figure 3C. The alternative designs were,
 341 crossing {environments} x {intercroops} x {N-fertilization}, {FR 9, FR 5, FR 22} x
 342 {dw/pea} x {N0, N} and {FR 22, FR 20, FR 16} x {dw/fb} x {N0, N}.

343 Discussion

344 One key reason to use agricultural data is to improve knowledge in crop science, as in other
 345 scientific fields. This can be generalized with the Data, Information, Knowledge and Wisdom
 346 pyramid (Ackoff, 1989), which describes the continuum between data and the knowledge
 347 it provides. Thus, the issue is to use appropriate methods based on the available data to
 348 provide insights and understanding of a studied system’s functioning. Depending on whether
 349 data come from experimental data or from scientific publications, methods related to global
 350 datasets or meta-analysis, respectively, will be used (Makowski et al., 2014), ~~and both.~~
 351 Both are useful for studying global issues in agronomy (Table 3). Two important issues arise
 352 from this observation: data availability and the knowledge that one wants to provide.

353 **Table 3. Overview of a comparison between meta-analysis and global datasets.**

| Criterion | Meta-analysis | Global datasets |
|--|---|---|
| Scope | All practices studied in multiple scientific publications | All practices tested in multiple experiments |
| Time required to collect and tidy the data | Long to very long (dozen to hundreds of hours) | Very long |
| Variables used | Often standard variables (e.g. <u>e.g.</u> yield, nitrogen fertilization) | All available observations (e.g. <u>e.g.</u> agronomic practices, phenotypic measurements, climate) |
| Number of observations | Moderate to large (dozens to hundreds) | Large (hundreds to thousands) |
| Reuse | Possible, but limited to the present variables | Possible once the data are formatted |
| Data sources | Scientific publications | Experimental files |

354 In meta-analysis, data are available because they are already published, even if it takes a
 355 long time to retrieve them. Conducting a meta-analysis is thus time-consuming, especially

356 the pre-analysis search and development of the database, which represent around 60% of the
357 working time (Allen and Olkin, 1999). Meta-analysis requires identifying and extracting the
358 values of interest from scientific publications, while being cautious to avoid potential bias.

359 In contrast, building global datasets requires interacting with the research teams that
360 conducted the experiments and adapting their raw experimental files to a standard format
361 (Figure 1). This step itself is very likely to necessitate more time than meta-analysis data
362 processing step—, and would greatly benefit from improved upstream data standardization
363 practices (Krajewski et al., 2015). The main advantage of global datasets in biology is
364 that they consist of phenotypic observations, which means that the studied processes are
365 potentially observed at lower levels than in meta-analysis. In this sense, global datasets
366 could enable further investigation of potential causalities based on correlations in the data
367 (Garside and Bell, 2011; Gunawardena, 2014). Additionally, since agronomic global datasets
368 contain plant-related variables measured at multiple organizational levels (~~e.g.~~e.g. organ,
369 plant, crop), they can target a wide audience for data reuse. For instance, researchers
370 developing functional–structural plant models (Louarn et al., 2020) may be interested in
371 variables measured at the plant scale (~~e.g.~~e.g. number of tillers, inter-node length, plant
372 height), while those who develop crop models to predict ~~yields~~yield (Berghuijs et al., 2021)
373 may be interested in variables measured at the crop scale (~~e.g.~~e.g. crop biomass, crop
374 height).

375 Alternatively, global datasets might have a role in increasing the discovery and use of
376 non-published experimental data. In our case, almost 30% of the experimental data gathered
377 have not been published through a research article. Bringing them together with other
378 experiments valued the time and energy required to conduct those field experiments. It
379 was also a friction point, since researchers may be reluctant to share unpublished data. For
380 instance, in our use case, 11 of the 37 experiments were not included in published articles or
381 database before this initiative, while each is now described within the global dataset (Gaudio
382 et al., 2023) and linked back groups leading field experiments in 1-4 scientific publications
383 (Gaudio et al., 2021; Louarn et al., 2021; Mahmoud et al., 2022; Meunier et al., 2022). Based
384 on the global dataset developed in this study, Gaudio et al. (2021) extracted a subset of 28
385 experiments to assess the influence of intercropping on the relation between plant biomass
386 and grain yield; Louarn et al. (2021) extracted a subset of 15 experiments to validate the
387 adaptation of Nitrogen Nutrition Index (NNI) to intercropping; Mahmoud et al. (2022)
388 extracted a subset of 11 experiments to assess the influence of ~~N~~nitrogen fertilization on
389 plant-plant interactions in intercrops; and Meunier et al. (2022) extracted a subset of 31
390 experiments to calibrate a statistical model used in a modeling chain to predict ecosystem
391 services as a function of the species associated in cereal-legume intercrops.

392 We argue that crop science can benefit from global datasets because they decrease the cost
393 of data (reuse) and increase the reproducibility of studies along with open data science
394 tools (Lowndes et al., 2017). Ultimately, global datasets contribute to new findings through
395 joint analysis of multiple experiments - a key consideration given the pressing need for
396 consolidating results in the context of an increasingly variable and changing climate. Despite
397 these needs for advancements, the challenges associated with the data standardization and
398 proprietary rights present significant obstacles to the ~~utilization~~-building of these global
399 datasets in crop science. A tighter integration between experimental and modeling research
400 communities is the first step in a way forward.

401 **Acknowledgements**

402 We thank the entire technical staff of the different research teams who shared their data,
403 for all the huge work they have done, without which this paper and the associated dataset
404 would not exist. We thank Michael and Michelle Corson for their helpful comments and
405 English revision, and the three reviewers for their valuable comments and corrections which
406 highly contribute to improve the manuscript.

407 **Funding**

408 This research was supported by the French National Research Agency under the Investments
409 for the Future Program (referred to as ANR-16-CONV-0004 and ANR-20-PCPA-0006) and
410 by the European Research Council under the European Union's Horizon Europe research
411 and innovation program in the framework of the IntercropValuES (Developing Intercropping
412 for agrifood Value chains and Ecosystem Services delivery in Europe and Southern countries,
413 <https://intercropvalues.eu/>) ~~project~~-starting from November 2022 [grant number 101081973].
414 ~~We thank Michael and Michelle Corson for their helpful comments and English revision.~~

415 **Conflict of interest disclosure**

416 The authors have no relevant financial or non-financial interests to disclose. On behalf of all
417 authors, the corresponding author states that there is no conflict of interest.

418 Author Contributions

419 All authors contributed to funding acquisition, data collection and formatting, writing and
420 editing the manuscript.

421 Data Availability

422 The global dataset is available on Zenodo open data repository (Gaudio et al., 2023).

423 References

- 424 Ackoff, R.L. 1989. From data to wisdom. *Journal of applied systems analysis* 16(1): 3–9.
- 425 Allen, I.E., and I. Olkin. 1999. Estimating Time to Conduct a Meta-analysis From Number
426 of Citations Retrieved. *JAMA* 282(7): 634–635. doi: [10.1001/JAMA.282.7.634](https://doi.org/10.1001/JAMA.282.7.634).
- 427 Barillot, R., D. Combes, S. Pineau, P. Huynh, and A.J. Escobar-Gutierrez. 2014. Comparison
428 of the morphogenesis of three genotypes of pea (*Pisum sativum*) grown in pure stands
429 and wheat-based intercrops. *Aob Plants* 6: plu006. doi: [https://doi.org/10.1093/aobpla/
430 plu006](https://doi.org/10.1093/aobpla/plu006).
- 431 Bedoussac, L., and E. Justes. 2010a. Dynamic analysis of competition and complementarity
432 for light and N use to understand the yield and the protein content of a durum wheat–
433 winter pea intercrop. *Plant and Soil* 330(1-2): 37–54. doi: [https://doi.org/10.1007/
434 s11104-010-0303-8](https://doi.org/10.1007/s11104-010-0303-8).
- 435 Bedoussac, L., and E. Justes. 2010b. The efficiency of a durum wheat-winter pea intercrop
436 to improve yield and wheat grain protein concentration depends on N availability during
437 early growth. *Plant and Soil* 330(1-2): 19–35. doi: [https://doi.org/10.1007/s11104-009-
438 0082-2](https://doi.org/10.1007/s11104-009-0082-2).
- 439 Berghuijs, H.N.C., M. Weih, W. Van Der Werf, A.J. Karley, E. Adam, et al. 2021. Calibrating
440 and testing APSIM for wheat-faba bean pure cultures and intercrops across Europe.
441 *Field Crops Research* 264: 108088. doi: [10.1016/j.fcr.2021.108088](https://doi.org/10.1016/j.fcr.2021.108088).
- 442 Broman, K.W., and K.H. Woo. 2018. Data organization in spreadsheets. *The American*
443 *Statistician* 72(1): 2–10. doi: [10.1080/00031305.2017.1375989](https://doi.org/10.1080/00031305.2017.1375989).
- 444 Bron, C., and J. Kerbosch. 1973. Algorithm 457: Finding All Cliques of an Undirected
445 Graph [H]. *Communications of the ACM* 16(9): 575–577. doi: [10.1145/362342.362367](https://doi.org/10.1145/362342.362367).
- 446 Casler, M.D. 2015. Fundamentals of experimental design: Guidelines for designing suc-
447 cessful experiments. *Agronomy Journal* 107(2): 692–705. doi: [https://doi.org/10.2134/
448 agronj2013.0114](https://doi.org/10.2134/agronj2013.0114).

- 449 Corre-Hellou, G., J. Fustec, and Y. Crozat. 2006. Interspecific competition for soil N and its
450 interaction with N₂ fixation, leaf expansion and crop growth in pea-barley intercrops.
451 *Plant and Soil* 282(1-2): 195–208. doi: <https://doi.org/10.1007/s11104-005-5777-4>.
- 452 Cruz, S.M.S. da, and J.A.P. do Nascimento. 2019. Towards integration of data-driven
453 agronomic experiments with data provenance. *Computers and Electronics in Agriculture*
454 161(September 2018): 14–28. doi: [10.1016/j.compag.2019.01.044](https://doi.org/10.1016/j.compag.2019.01.044).
- 455 Duru, M., O. Therond, G. Martin, R. Martin-Clouaire, M.-A. Magne, et al. 2015. How
456 to implement biodiversity-based agriculture to enhance ecosystem services: A review.
457 *Agronomy for Sustainable Development* 35(4): 1259–1281. doi: [10.1007/s13593-015-0306-](https://doi.org/10.1007/s13593-015-0306-1)
458 [1](https://doi.org/10.1007/s13593-015-0306-1).
- 459 Garside, A.L., and M.J. Bell. 2011. Growth and yield responses to amendments to the
460 sugarcane monoculture: Towards identifying the reasons behind the response to breaks.
461 *Crop and Pasture Science* 62(9): 776–789. doi: [10.1071/CP11055](https://doi.org/10.1071/CP11055).
- 462 Gaudio, N., R. Mahmoud, L. Bedoussac, E. Justes, E.-P. Journet, et al. 2023. A global
463 dataset gathering 37 field experiments involving cereal-legume intercrops and their
464 corresponding sole crops. doi: [10.5281/zenodo.8081577](https://doi.org/10.5281/zenodo.8081577).
- 465 Gaudio, N., C. Violle, X. Gendre, F. Fort, R. Mahmoud, et al. 2021. Interspecific interactions
466 regulate plant reproductive allometry in cereal–legume intercropping systems. *Journal of*
467 *Applied Ecology* 58(11): 2579–2589. doi: <https://doi.org/10.1111/1365-2664.13979>.
- 468 Glass, G.V. 1976. Primary, secondary, and meta-analysis of research. *Educational researcher*
469 5(10): 3–8.
- 470 Gunawardena, J. 2014. Models in biology: 'Accurate descriptions of our pathetic thinking'.
471 *BMC Biology* 12(1): 1–11. doi: [10.1186/1741-7007-12-29/FIGURES/3](https://doi.org/10.1186/1741-7007-12-29/FIGURES/3).
- 472 Hauggaard-Nielsen, H., M. Gooding, P. Ambus, G. Corre-Hellou, Y. Crozat, et al. 2009a.
473 Pea-barley intercropping for efficient symbiotic N₂-fixation, soil N acquisition and use
474 of other nutrients in European organic cropping systems. *Field Crops Research* 113(1):
475 64–71. doi: <https://doi.org/10.1016/j.fcr.2009.04.009>.
- 476 Hauggaard-Nielsen, H., M. Gooding, P. Ambus, G. Corre-Hellou, Y. Crozat, et al. 2009b.
477 Pea-barley intercropping and short-term subsequent crop effects across European organic
478 cropping conditions. *Nutrient Cycling in Agroecosystems* 85(2): 141–155. doi: <https://doi.org/10.1007/s10705-009-9254-y>.
- 480 Hauggaard-Nielsen, H., B. Jørnsgaard, J. Kinane, and E.S. Jensen. 2008. Grain legume–
481 cereal intercropping: The practical application of diversity, competition and facilitation
482 in arable and organic cropping systems. *Renewable Agriculture and Food Systems* 23(1):
483 3–12. doi: <https://doi.org/10.1017/S1742170507002025>.
- 484 Jenkins, G.B., A.P. Beckerman, C. Bellard, A. Benitez-López, A.M. Ellison, et al. 2023.
485 Reproducibility in ecology and evolution: Minimum standards for data and code. *Ecology*
486 *and Evolution* 13(5). doi: [10.1002/ece3.9961](https://doi.org/10.1002/ece3.9961).

487 Kammoun, B., E.-P. Journet, E. Justes, and L. Bedoussac. 2021. Cultivar Grain Yield in
488 Durum Wheat-Grain Legume Intercrops Could Be Estimated From Sole Crop Yields
489 and Interspecific Interaction Index. *Frontiers in Plant Science* 12: 2191. doi: <https://doi.org/10.3389/fpls.2021.733705>.
490

491 Kattge, J., S. Diaz, S. Lavorel, I.C. Prentice, P. Leadley, et al. 2011. TRY—a global database
492 of plant traits. *Global change biology* 17(9): 2905–2935.

493 Knudsen, M.T., H. Hauggaard-Nielsen, B. Jørgensen, and E.S. Jensen. 2004. Comparison
494 of interspecific competition and N use in pea-barley, faba bean-barley and lupin-barley
495 intercrops grown at two temperate locations. *Journal of Agricultural Science* 142: 617–627.
496 doi: <https://doi.org/10.1017/S0021859604004745>.

497 [Krajewski, P., D. Chen, H. Ćwiek, A.D.J. van Dijk, F. Fiorani, et al. 2015. Towards
498 recommendations for metadata and data handling in plant phenotyping. *Journal of
499 Experimental Botany* 66\(18\): 5417–5427. doi: <https://doi.org/10.1093/jxb/erv271>.](https://doi.org/10.1093/jxb/erv271)

500 [Launay, M., N. Brisson, S. Satger, H. Hauggaard-Nielsen, G. Corre-Hellou, et al. 2009.
501 pre Exploring options for managing strategies for pea-barley intercropping using a modeling
502 approach. *European Journal of Agronomy* 31\(2\): 85–98. doi: \[https://doi.org/10.1016/j.
503 eja.2009.04.002\]\(https://doi.org/10.1016/j.eja.2009.04.002\).](https://doi.org/10.1016/j.eja.2009.04.002)

504 Lawler, E.L., J.K. Lenstra, and A.H.G. Rinnooy Kan. 1980. Generating All Maximal
505 Independent Sets: NP-Hardness and Polynomial-Time Algorithms. *SIAM Journal on
506 Computing* 9(3): 558–565. doi: 10.1137/0209042.

507 Licker, R., M. Johnston, J.A. Foley, C. Barford, C.J. Kucharik, et al. 2010. Mind the gap:
508 How do climate and agricultural management explain the 'yield gap' of croplands around
509 the world? *Global Ecology and Biogeography* 19(6): 769–782. doi: 10.1111/j.1466-
510 8238.2010.00563.x.

511 Lobell, D.B., J.M. Deines, and S.D. Tommaso. 2020. Changes in the drought sensitivity of
512 US maize yields. *Nature Food* 1(11): 729–735. doi: 10.1038/s43016-020-00165-w.

513 Louarn, G., R. Barillot, Di. Combes, and A. Escobar-Gutiérrez. 2020. Towards intercrop
514 ideotypes: Non-random trait assembly can promote overyielding and stability of species
515 proportion in simulated legume-based mixtures. *Annals of Botany* 126(4): 671–685. doi:
516 10.1093/aob/mcaa014.

517 Louarn, G., L. Bedoussac, N. Gaudio, E.P. Journet, D. Moreau, et al. 2021. Plant nitrogen
518 nutrition status in intercrops— a review of concepts and methods. *European Journal of
519 Agronomy* 124: 126229. doi: 10.1016/J.EJA.2021.126229.

520 Lowndes, J.S.S., B.D. Best, C. Scarborough, J.C. Afflerbach, M.R. Frazier, et al. 2017. Our
521 path to better science in less time using open data science tools. *Nature Ecology &
522 Evolution* 1(6): 1–7. doi: <https://doi.org/10.1038/s41559-017-0160>.

523 Maat, H. 2011. The history and future of agricultural experiments. *NJAS - Wageningen*

524 Journal of Life Sciences 57(3-4): 187–195. doi: 10.1016/j.njas.2010.11.001.

525 Mahmoud, R., P. Casadebaig, N. Hilgert, L. Alletto, G.T. Freschet, et al. 2022. Species choice
526 and n fertilization influence yield gains through complementarity and selection effects in
527 cereal-legume intercrops. *Agronomy for sustainable development*. doi: 10.1007/s13593-
528 022-00754-y.

529 Makowski, D., T. Nesme, F. Papy, and T. Doré. 2014. Global agronomy, a new field
530 of research. A review. *Agronomy for Sustainable Development* 34(2): 293–307. doi:
531 10.1007/s13593-013-0179-0.

532 Meunier, C., L. Alletto, L. Bedoussac, J.E. Bergez, P. Casadebaig, et al. 2022. A modelling
533 chain combining soft and hard models to assess a bundle of ecosystem services provided
534 by a diversity of cereal-legume intercrops. *European Journal of Agronomy* 132(October
535 2021). doi: 10.1016/j.eja.2021.126412.

536 Naudin, C., G. Corre-Hellou, S. Pineau, Y. Crozat, and M.-H. Jeuffroy. 2010. The effect
537 of various dynamics of N availability on winter pea-wheat intercrops: Crop growth,
538 N partitioning and symbiotic N-2 fixation. *Field Crops Research* 119(1): 2–11. doi:
539 <https://doi.org/10.1016/j.fcr.2010.06.002>.

540 Naudin, C., H.M.G. van der Werf, M.-H. Jeuffroy, and G. Corre-Hellou. 2014. Life cycle
541 assessment applied to pea-wheat intercrops: A new method for handling the impacts of
542 co-products. *Journal of Cleaner Production* 73: 80–87. doi: <https://doi.org/10.1016/j.jclepro.2013.12.029>.

544 Newman, S.J., and R.T. Furbank. 2021. A multiple species, continent-wide, million-
545 phenotype agronomic plant dataset. *Scientific Data* 8(1): 1–8. doi: 10.1038/s41597-021-
546 00898-8.

547 Pelzer, E., M. Bazot, L. Guichard, and M.-H. Jeuffroy. 2016. Crop Management Affects the
548 Performance of a Winter Pea–Wheat Intercrop. *Agronomy Journal* 108(3): 1089–1100.
549 doi: <https://doi.org/10.2134/agronj2015.0440>.

550 Phillips, C.A., K. Wang, E.J. Baker, J.A. Bubier, E.J. Chesler, et al. 2019. On Finding and
551 enumerating maximal and maximum k-partite cliques in k-partite graphs. *Algorithms*
552 12(1). doi: 10.3390/a12010023.

553 Popkin, G. 2019. Data sharing and how it can benefit your scientific career. *Nature* 569(7756):
554 445–447. doi: 10.1038/d41586-019-01506-x.

555 Senft, M., U. Stahl, and N. Svoboda. 2022. Research data management in agricultural
556 sciences in germany: We are not yet where we want to be (C. Pulvento, editor). *PLOS*
557 *ONE* 17(9): e0274677. doi: 10.1371/journal.pone.0274677.

558 Sparks, A.H. 2018. Nasapower: A NASA POWER Global Meteorology, Surface Solar
559 Energy and Climatology Data Client for R. *Journal of Open Source Software* 3(30):
560 1035. doi: 10.21105/joss.01035.

Tang, X., S.A. Placella, F. Dayde, L. Bernard, A. Robin, et al. 2016. Phosphorus availability

561 [pre](#)

562 and microbial community in the rhizosphere of intercropped cereal and legume along a P-

563 fertilizer gradient. *Plant and Soil* 407(1-2): 119–134. doi: [https://doi.org/10.1007/s11104-](https://doi.org/10.1007/s11104-016-2949-3)

564 [016-2949-3](https://doi.org/10.1007/s11104-016-2949-3).

565 Tardieu, F. 2020. Educated big data to study sensitivity to drought. *Nature Food* 1(11):

566 669–670. doi: [10.1038/s43016-020-00187-4](https://doi.org/10.1038/s43016-020-00187-4).

567 Viguier, L., L. Bedoussac, E.-P. Journet, and E. Justes. 2018. Yield gap analysis extended

568 to marketable grain reveals the profitability of organic lentil-spring wheat intercrops.

569 *Agronomy for Sustainable Development* 38(4): 39. doi: [https://doi.org/10.1007/s13593-](https://doi.org/10.1007/s13593-018-0515-5)

570 [018-0515-5](https://doi.org/10.1007/s13593-018-0515-5).

571 White, J.W., and F.K. Van Evert. 2008. Publishing agronomic data. *Agronomy Journal*

572 100(5): 1396–1400. doi: [10.2134/agronj2008.0080F](https://doi.org/10.2134/agronj2008.0080F).

573 Wickham, H. 2014. Tidy data. *Journal of Statistical Software* 59(10). doi:

574 [10.18637/jss.v059.i10](https://doi.org/10.18637/jss.v059.i10).

575 Wickham, H., and G. Grolemund. 2016. R for data science: Import, tidy, transform,

576 visualize, and model data. " O'Reilly Media, Inc."

577 Wilkinson, M.D., M. Dumontier, Ij.J. Aalbersberg, G. Appleton, M. Axton, et al. 2016.

578 Comment: The FAIR Guiding Principles for scientific data management and stewardship.

579 *Scientific Data* 3: 1–9. doi: [10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).

580 Wilson, G., J. Bryan, K. Cranston, J. Kitzes, L. Nederbragt, et al. 2017. Good enough

581 practices in scientific computing. *PLOS Computational Biology* 13(6): e1005510. doi:

582 [10.1371/JOURNAL.PCBI.1005510](https://doi.org/10.1371/JOURNAL.PCBI.1005510).

583 Zamir, D. 2013. Where Have All the Crop Phenotypes Gone? *PLoS Biology* 11(6): 1–4. doi:

584 [10.1371/journal.pbio.1001595](https://doi.org/10.1371/journal.pbio.1001595).