



HAL
open science

A workflow for processing global datasets: application to intercropping

Rémi Mahmoud, Pierre Casadebaig, Nadine Hilgert, Noémie Gaudio

► To cite this version:

Rémi Mahmoud, Pierre Casadebaig, Nadine Hilgert, Noémie Gaudio. A workflow for processing global datasets: application to intercropping. *Peer Community Journal*, 2024, 4, pp.e24. 10.24072/pcjournal.389 . hal-04145269v4

HAL Id: hal-04145269

<https://hal.inrae.fr/hal-04145269v4>

Submitted on 27 Mar 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



Peer Community Journal

Section: Mathematical & Computational Biology

Research article

Published
2024-02-29

Cite as

Rémi Mahmoud, Pierre Casadebaig, Nadine Hilgert and Noémie Gaudio (2024) A workflow for processing global datasets: application to intercropping, Peer Community Journal, 4: e24.

Correspondence
pierre.casadebaig@inrae.fr

Peer-review

Peer reviewed and recommended by PCI Mathematical & Computational Biology,
<https://doi.org/10.24072/pci.mcb.100197>



This article is licensed under the Creative Commons Attribution 4.0 License.

A workflow for processing global datasets: application to intercropping

Rémi Mahmoud^{,1}, Pierre Casadebaig^{,1}, Nadine Hilgert², and Noémie Gaudio^{,1}

Volume 4 (2024), article e24

<https://doi.org/10.24072/pcjournal.389>

Abstract

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

¹AGIR, Univ. Toulouse, INRAE, Castanet-Tolosan, France, ²MISTEA, Univ. Montpellier, INRAE, Institut Agro, Montpellier, France



Introduction

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

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

However, global datasets come with their own challenges. Assembling these datasets requires extensive data collection, standardization, and homogenization across diverse experiments conducted by different research teams (Makowski et al., 2014; White et al., 2008). This tedious curation step is an undervalued task, whose duration could be reduced from the adoption of good practices upstream. Recent efforts and international initiatives aimed at opening and standardizing data are emerging, highlighting that data standardization is crucial for improving the interpretation of experimental results and the generalization of knowledge acquisition. It also facilitates statistical meta-analysis and data publication (Krajewski et al., 2015). However, datasets for plant and crop measurements in controlled field trials are still scarce in public databases. The different field experiments gathered often have diverse objectives, leading to unbalanced and incomplete designs. Confounding factors, *i.e.* the unintended mixing of two or more effects making them indistinguishable, can also be challenging (Casler, 2015). Consequently, using and analyzing global datasets require a thorough understanding of the dataset, judicious interpretation of the results, identification of balanced data subsets for specific re-search questions, and acceptance that the effects of some factors may remain indistinguishable. Therefore, the application of statistical learning techniques on global datasets is only feasible after extensive data pre-processing.

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

To achieve this, we recommend adopting practices for designing and analyzing global datasets that align with tidy data (Broman et al., 2018; Wickham, 2014) and FAIR principles (Findable, Accessible, Interoperable, and Reusable) (Wilkinson et al., 2016). As a use case, we illustrate the design of a global dataset for intercropping systems, in which at least two crop species are

grown in the same field for a significant part of their growth cycle. We describe the main steps involved in designing a global dataset gathering 37 intercropping experiments across Europe. We also describe and apply an original method to identify complete factorial design subsets of interest. This methodological development was aimed at helping the potential collaborators to explore and get an overview of the dataset as a function of their factor of interest, a key step in assisting further modeling and analysis steps.

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

Design steps of global datasets

This section presents the generic steps involved in designing a global dataset. As the gathering, cleaning, and formatting of the spare source datasets is time-consuming, we followed tidy data specifications (Wickham, 2014) and a global data science workflow as presented by Wickham and Grolemund (2016) (Figure 1).

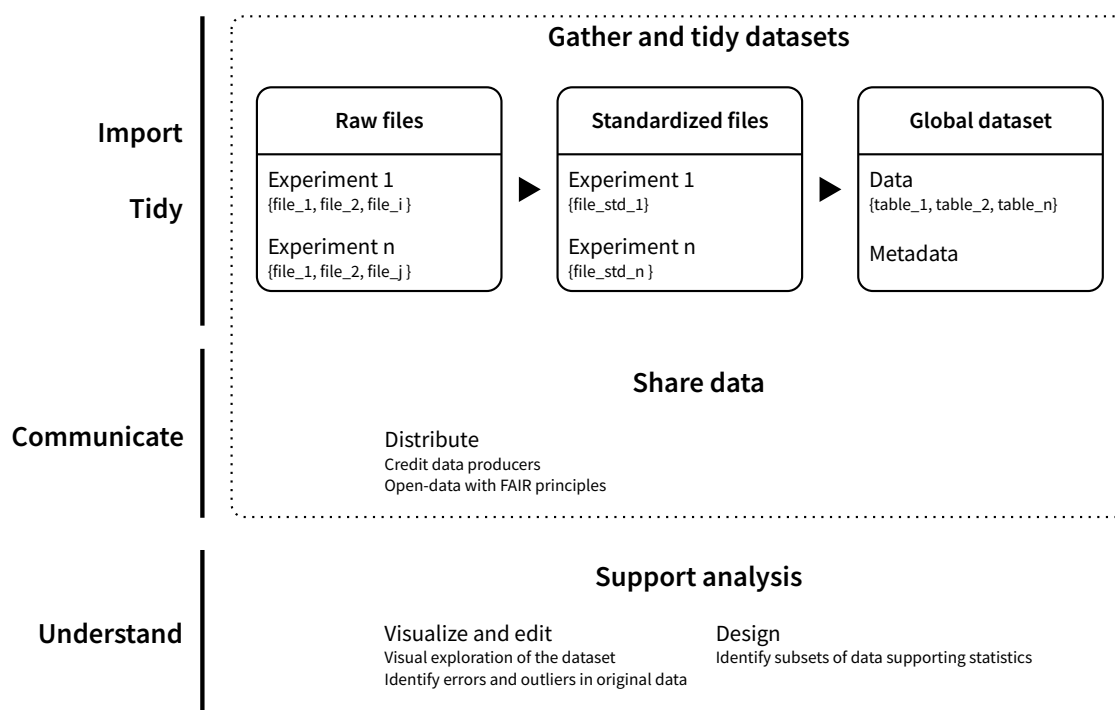


Figure 1 - Main steps for designing global datasets. The left column corresponds to a classical data science workflow. We adapted these steps for global dataset design specificities, to illustrate the importance of data gathering, tidying, and sharing (dotted frame). While some actions supporting subsequent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

1. Gather and tidy source datasets

1.1. Conceptual framework

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

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

kin, 2019).

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

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

1.2. Case study

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

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

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

Finally, all the files were pooled together using R software, to create one global table per data category, *i.e.* four tables related respectively to climate, crop measurements, agricultural prac-

tices and global information describing the site (Figure 2). Overall, the global dataset contained 308 and 299 statistical individuals (defined as a unique combination of {site * year * management}) in intercrop and sole crop, respectively (Table 1). The number of plant characteristics was much larger (33351 observations, among which 12896 were measured in sole crops and 20455 in intercrops), since several variables were measured at the crop scale, sometimes several times during the crop cycle.

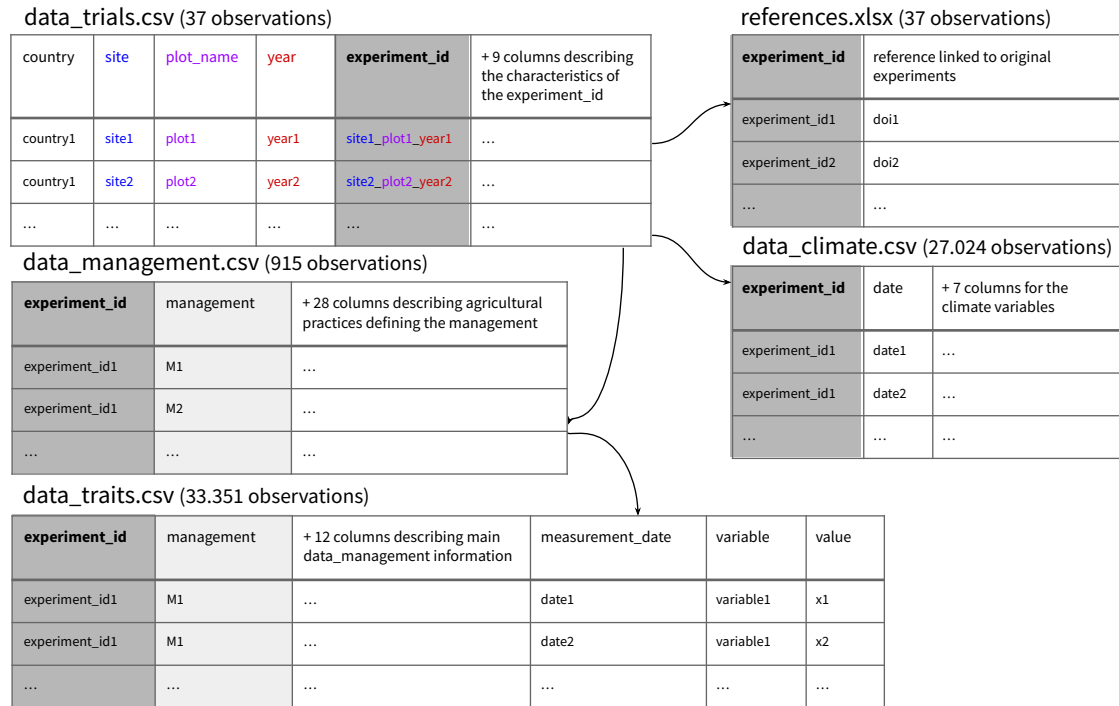


Figure 2 - Representation of the relationships between tables identified in the global dataset. Five tables were defined to organize data, all sharing a common identifier (*experiment_id*, which is the concatenation of the *site_plot_year* of each experiment). The table *data_trials.csv* provides the 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).

Table 1 - Overview of the diversity of the treatments in the global dataset by factors (columns) and experiments (rows). Within each column, each colored rectangle is a level of the factor considered. For instance, the two colors for the *Mixing pattern* indicate that the two species intercropped were sown in alternate rows or within the row; the two colors for the *Nitrogen (N) fertilization* indicate that the experiment included at least two N-treatments (no fertilization and N-fertilization, the latter of which may include several amounts of N); regarding *Species mixture*, the number of colors indicates the number of different species mixtures included in a given experiment. A rectangle in a given row and column indicates that the corresponding 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]	[Cyan]
Jynde vad_jyn_2002	24	12	[Red]	[Green]	[Cyan]
Jynde vad_jyn_2001	24	12	[Red]	[Green]	[Cyan]
Grignon_inrae_2017	19	6	[Red]	[Green]	[Cyan]
Grignon_inrae_2010	16	8	[Red]	[Green]	[Cyan]
Grignon_inrae_2009	15	8	[Red]	[Green]	[Cyan]
Grignon_inrae_2008	27	5	[Red]	[Green]	[Cyan]
Grignon_inrae_2007	30	7	[Red]	[Green]	[Cyan]
Copenhagen_hbg_2003	24	10	[Red]	[Green]	[Cyan]
Copenhagen_hbg_2002	24	11	[Red]	[Green]	[Cyan]
Copenhagen_hbg_2001	24	12	[Red]	[Green]	[Cyan]
Auz_ZN_2012	58	24	[Red]	[Green]	[Cyan]
Auz_TO_2016	86	18	[Red]	[Green]	[Cyan]
Auz_TO_2013	93	24	[Red]	[Green]	[Cyan]
Auz_TE_2006	13	20	[Red]	[Green]	[Cyan]
Auz_SGs_2007	66	23	[Red]	[Green]	[Cyan]
Auz_PP_2011	20	20	[Red]	[Green]	[Cyan]
Auz_pk_2011	18	18	[Red]	[Green]	[Cyan]
Auz_marinette_2_2015	85	13	[Red]	[Green]	[Cyan]
Auz_marinette_1_2015	22	13	[Red]	[Green]	[Cyan]
Auz_ochard_2010	60	21	[Red]	[Green]	[Cyan]
Angers_thorigne_2009	11	12	[Red]	[Green]	[Cyan]
Angers_thorigne_2008	15	14	[Red]	[Green]	[Cyan]
Angers_thorigne_2007	11	12	[Red]	[Green]	[Cyan]
Angers_thorigne_2006	6	8	[Red]	[Green]	[Cyan]
Angers_thorigne_2004	6	10	[Red]	[Green]	[Cyan]
Angers_thorigne_2003	6	10	[Red]	[Green]	[Cyan]
Angers_jailliere_2008	22	16	[Red]	[Green]	[Cyan]
Angers_jailliere_2007	14	16	[Red]	[Green]	[Cyan]
Angers_fnams_2003	12	10	[Red]	[Green]	[Cyan]
Angers_fnams_2002	4	8	[Red]	[Green]	[Cyan]
Angers_brainsurlauthion_2011	10	5	[Red]	[Green]	[Cyan]

2. Share organized data

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

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

This global dataset, as well as the metadata associated, are available on a data repository in a FAIR way (Gaudio, Mahmoud, et al., 2023). Out of the 37 experiments gathered, 11 have never been valued before. Additional details on experimental designs and management practices are reported in the reference publications for 26 of the 37 experiments (Barillot et al., 2014; Bedoussac et al., 2010a,b; Corre-Hellou et al., 2006; Hauggaard-Nielsen, Gooding, Am-

bus, Corre-Hellou, Crozat, Dahlmann, Dahlmann, et al., 2009; Hauggaard-Nielsen, Gooding, Ambus, Corre-Hellou, Crozat, Dahlmann, Dibet, et al., 2009; Hauggaard-Nielsen, Jørnsgaard, et al., 2008; Kammoun et al., 2021; Knudsen et al., 2004; Launay et al., 2009; Naudin, Corre-Hellou, et al., 2010; Naudin, Werf, et al., 2014; Pelzer et al., 2016; Tang et al., 2016; Viguier et al., 2018).

3. Support new analysis

3.1. Conceptual framework

Once the data are in a tractable format, visual exploration allows for a comprehensive overview of data patterns, aiding in the identification of anomalies such as errors and outliers that may not be immediately apparent through numerical analysis alone. Later, additional processes are required to render the dataset operational for analytical and modeling studies, such as data imputation, dimension reduction, or data normalization. Because these steps depend largely on the chosen analytical workflow, they are not directly included in the communicated open datasets, but rather tailored by the subsequent analytical team (Figure 1). Nonetheless, sharing methods can support the future reuse of the dataset. In our case in crop ecology, we illustrated this step with the development of an original method aiming at identifying subsets in the overall dataset corresponding to complete factorial designs.

3.2. Case study

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

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

Table 2 - Definitions in graph theory used in the present study.

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

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

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

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

Application Here, we illustrate this method with two datasets : (1) a theoretical one, where we generated an unbalanced design of five environments, five intercrops, and two management levels (Figure 3A); and (2) a practical one, corresponding to the global dataset presented in this study (Figure 3B and 3C).

When applied on the theoretical unbalanced design (Figure 3A), this method identified 8 maximal 3-partite cliques, each of these designs having different number of modalities in considered factors (environment, intercrops or management). There is only one design maximizing the number of environments, and no factorial design was found with two levels per factor.

We considered two examples for the application on the agronomic global dataset. In the first one, we searched for any number of intercrops observed at least in two environments. Two designs were identified: the one with the most environmental modalities is illustrated in Figure 3B; the alternative design was, crossing {environments} x {intercrops}, {FR_22, FR_21} x {dw/pea, dw/fb}. The second example was the same request with an additional constraint on two levels of nitrogen (N) fertilization. In this case, three designs were identified, the largest one being illustrated in Figure 3C. The alternative designs were, crossing {environments} x {intercrops} x {N-fertilization}, {FR_9, FR_5, FR_22} x {dw/pea} x {N0, N} and {FR_22, FR_20, FR_16} x {dw/fb} x {N0, N}.

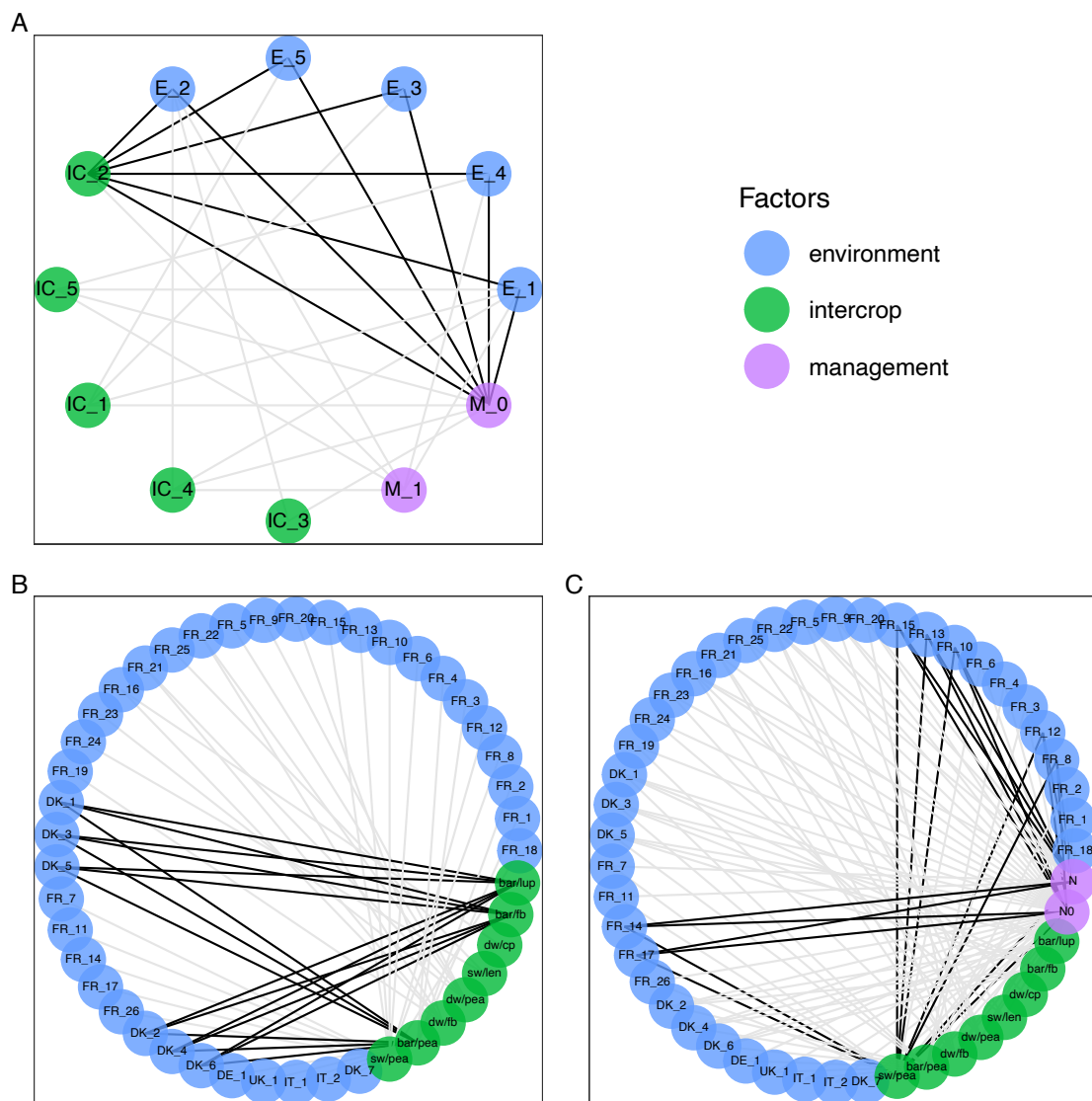


Figure 3 - Three maximal k -cliques that represent distinct complete factorial designs within theoretical (A) and experimental (B-C) unbalanced designs. Black edges represent the edges of the cliques and gray edges represent the factor combinations appearing in the initial design. In the case A, we generated a random unbalanced design for three factors and illustrated the 3-clique maximizing the number of environments. The experimental design in the cases B and C corresponds to the aggregation of the 37 experimentations (blue nodes). In case B, we searched for any intercrop observed at least in two environments. In case C, there was an additional constraint on two levels of nitrogen (N) fertilization. Countries were abbreviated with their ISO 3166 codes; species were abbreviated as barley (*bar*), chickpea (*cp*), durum wheat (*dw*), faba bean (*fb*), lentil (*len*), lupin (*lup*), soft wheat (*sw*); nitrogen fertilization was abbreviated as NO for no fertilization, and N for fertilization.

Discussion

One key reason to use agricultural data is to improve knowledge in crop science, as in other scientific fields. This can be generalized with the Data, Information, Knowledge and Wisdom pyramid (Ackoff, 1989), which describes the continuum between data and the knowledge it provides. Thus, the issue is to use appropriate methods based on the available data to provide insights and understanding of a studied system's functioning. Depending on whether data come from experimental data or from scientific publications, methods related to global datasets or meta-analysis,

respectively, will be used (Makowski et al., 2014). Both are useful for studying global issues in agronomy (Table 3). Two important issues arise from this observation: data availability and the knowledge that one wants to provide.

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. yield, nitrogen fertilization)	All available observations (e.g. 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

In meta-analysis, data are available because they are already published, even if it takes a long time to retrieve them. Conducting a meta-analysis is thus time-consuming, especially the pre-analysis search and development of the database, which represent around 60% of the working time (Allen et al., 1999). Meta-analysis requires identifying and extracting the values of interest from scientific publications, while being cautious to avoid potential bias.

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

Alternatively, global datasets might have a role in increasing the discovery and use of non-published experimental data. In our case, almost 30% of the experimental data gathered have not been published through a research article. Bringing them together with other experiments valued the time and energy required to conduct those field experiments. It was also a friction point, since researchers may be reluctant to share unpublished data. For instance, in our use case, 11 of the 37 experiments were not included in published articles or database before this initiative, while each is now described within the global dataset (Gaudio, Mahmoud, et al., 2023) and linked back groups leading field experiments in 1-4 scientific publications (Gaudio, Violle, et al., 2021; Louarn, Bedoussac, et al., 2021; Mahmoud et al., 2022; Meunier et al., 2022). Based on the global dataset developed in this study, Gaudio, Violle, et al. (2021) extracted a subset of 28 experiments to assess the influence of intercropping on the relation between plant biomass and grain yield; Louarn, Bedoussac, et al. (2021) extracted a subset of 15 experiments to validate the adaptation of Nitrogen Nutrition Index (NNI) to intercropping; Mahmoud et al. (2022) extracted a subset of 11 experiments to assess the influence of nitrogen fertilization on plant-plant interactions in intercrops; and Meunier et al. (2022) extracted a subset of 31 experiments to calibrate

a statistical model used in a modeling chain to predict ecosystem services as a function of the species associated in cereal-legume intercrops.

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

Acknowledgements

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

Preprint version 2 (<https://hal.science/hal-04145269>) of this article has been peer-reviewed and recommended by Peer Community In Mathematical and Computational Biology (<https://doi.org/10.24072/pci.mcb.100197>; Tannier, 2024).

Conflict of interest disclosure

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

Author Contributions

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

Data Availability

The dataset presented in this work is available on the Zenodo open data repository (<https://doi.org/10.5281/zenodo.8081577>; Gaudio, Mahmoud, et al., 2023).

Funding

This research was supported by the French National Research Agency under the Investments for the Future Program (referred to as ANR-16-CONV-0004 and ANR-20-PCPA-0006) and by the European Research Council under the European Union's Horizon Europe research and innovation program in the framework of the IntercropValuES (Developing Intercropping for agri-food Value chains and Ecosystem Services delivery in Europe and Southern countries, <https://intercropvalues.eu/>) starting from November 2022 [grant number 101081973].

References

- Ackoff RL (1989). From data to wisdom. *Journal of applied systems analysis* 16, 3–9.
- Allen IE and I Olkin (1999). Estimating Time to Conduct a Meta-analysis From Number of Citations Retrieved. *JAMA* 282, 634–635. <https://doi.org/10.1001/JAMA.282.7.634>.

- Barillot R, D Combes, S Pineau, P Huynh, and AJ Escobar-Gutierrez (2014). Comparison of the morphogenesis of three genotypes of pea (*Pisum sativum*) grown in pure stands and wheat-based intercrops. *Aob Plants* 6, plu006. <https://doi.org/10.1093/aobpla/plu006>.
- Bedoussac L and E Justes (2010a). Dynamic analysis of competition and complementarity for light and N use to understand the yield and the protein content of a durum wheat–winter pea intercrop. en. *Plant and Soil* 330, 37–54. <https://doi.org/10.1007/s11104-010-0303-8>.
- (2010b). The efficiency of a durum wheat–winter pea intercrop to improve yield and wheat grain protein concentration depends on N availability during early growth. en. *Plant and Soil* 330, 19–35. <https://doi.org/10.1007/s11104-009-0082-2>.
- Berghuijs HNC, M Weih, W Van Der Werf, AJ Karley, E Adam, AM Villegas-Fernández, LP Ki-aer, AC Newton, C Scherber, S Tavoletti, and G Vico (2021). Calibrating and testing APSIM for wheat-faba bean pure cultures and intercrops across Europe. *Field Crops Research* 264, 108088. <https://doi.org/10.1016/j.fcr.2021.108088>.
- Broman KW and KH Woo (2018). Data Organization in Spreadsheets. *The American Statistician* 72, 2–10. <https://doi.org/10.1080/00031305.2017.1375989>.
- Bron C and J Kerbosch (1973). Algorithm 457: Finding All Cliques of an Undirected Graph [H]. *Communications of the ACM* 16, 575–577. <https://doi.org/10.1145/362342.362367>.
- Casler MD (2015). Fundamentals of Experimental Design: Guidelines for Designing Successful Experiments. *Agronomy Journal* 107, 692–705. <https://doi.org/10.2134/agronj2013.0114>.
- Corre-Hellou G, J Fustec, and Y Crozat (2006). Interspecific competition for soil N and its interaction with N-2 fixation, leaf expansion and crop growth in pea–barley intercrops. *Plant and Soil* 282, 195–208. <https://doi.org/10.1007/s11104-005-5777-4>.
- Cruz SMSd and JAPd Nascimento (2019). Towards integration of data-driven agronomic experiments with data provenance. *Computers and Electronics in Agriculture* 161, 14–28. <https://doi.org/10.1016/j.compag.2019.01.044>.
- Duru M, O Therond, G Martin, R Martin-Clouaire, MA Magne, E Justes, EP Journet, JN Aubertot, S Savary, JE Bergez, and JP Sarthou (2015). How to implement biodiversity-based agriculture to enhance ecosystem services: a review. *Agronomy for Sustainable Development* 35, 1259–1281. <https://doi.org/10.1007/s13593-015-0306-1>.
- Garside AL and MJ Bell (2011). Growth and yield responses to amendments to the sugarcane monoculture: Towards identifying the reasons behind the response to breaks. *Crop and Pasture Science* 62, 776–789. <https://doi.org/10.1071/CP11055>.
- Gaudio N, R Mahmoud, L Bedoussac, E Justes, EP Journet, C Naudin, H Hauggaard-Nielsen, ES Jensen, E Pelzer, G Corre-Hellou, B Kammoun, L Viguier, R Barillot, A Couëdel, P Hinsinger, and P Casadebaig (2023). A global dataset gathering 37 field experiments involving cereal-legume intercrops and their corresponding sole crops. eng. <https://doi.org/10.5281/zenodo.8081577>.
- Gaudio N, C Violle, X Gendre, F Fort, R Mahmoud, E Pelzer, S Médiène, H Hauggaard-Nielsen, L Bedoussac, C Bonnet, G Corre-Hellou, A Couëdel, P Hinsinger, ES Jensen, EP Journet, E Justes, B Kammoun, I Litrico, N Moutier, C Naudin, and P Casadebaig (2021). Interspecific interactions regulate plant reproductive allometry in cereal–legume intercropping systems. en. *Journal of Applied Ecology* 58, 2579–2589. <https://doi.org/10.1111/1365-2664.13979>.
- Glass GV (1976). Primary, secondary, and meta-analysis of research. *Educational researcher* 5, 3–8.
- Gunawardena J (2014). Models in biology: 'Accurate descriptions of our pathetic thinking'. *BMC Biology* 12, 1–11. <https://doi.org/10.1186/1741-7007-12-29>.
- Hauggaard-Nielsen H, M Gooding, P Ambus, G Corre-Hellou, Y Crozat, C Dahlmann, C Dahlmann, Pv Fragstein, A Pristeri, M Monti, and ES Jensen (2009). Pea–barley intercropping for efficient symbiotic N-2-fixation, soil N acquisition and use of other nutrients in European organic cropping systems. *Field Crops Research* 113, 64–71. <https://doi.org/10.1016/j.fcr.2009.04.009>.
- Hauggaard-Nielsen H, M Gooding, P Ambus, G Corre-Hellou, Y Crozat, C Dahlmann, A Dibet, P von Fragstein, A Pristeri, M Monti, and ES Jensen (2009). Pea–barley intercropping and

- short-term subsequent crop effects across European organic cropping conditions. *Nutrient Cycling in Agroecosystems* 85, 141–155. <https://doi.org/10.1007/s10705-009-9254-y>.
- Hauggaard-Nielsen H, B Jørnsgaard, J Kinane, and ES Jensen (2008). Grain legume–cereal intercropping: The practical application of diversity, competition and facilitation in arable and organic cropping systems. *Renewable Agriculture and Food Systems* 23, 3–12. <https://doi.org/10.1017/S1742170507002025>.
- Jenkins GB, AP Beckerman, C Bellard, A Benítez-López, AM Ellison, CG Foote, AL Hufton, MA Lashley, CJ Lortie, Z Ma, AJ Moore, SR Narum, J Nilsson, B O’Boyle, DB Provete, O Razgour, L Rieseberg, C Riginos, L Santini, B Sibbett, and PR Peres-Neto (2023). Reproducibility in ecology and evolution: Minimum standards for data and code. *Ecology and Evolution* 13. <https://doi.org/10.1002/ece3.9961>.
- Kammoun B, EP Journet, E Justes, and L Bedoussac (2021). Cultivar Grain Yield in Durum Wheat–Grain Legume Intercrops Could Be Estimated From Sole Crop Yields and Interspecific Interaction Index. *Frontiers in Plant Science* 12, 2191. <https://doi.org/10.3389/fpls.2021.733705>.
- KATTGE J et al. (2011). TRY – a global database of plant traits. *Global Change Biology* 17, 2905–2935. <https://doi.org/10.1111/j.1365-2486.2011.02451.x>.
- Knudsen MT, H Hauggaard-Nielsen, B Jørnsgaard, and ES Jensen (2004). Comparison of interspecific competition and N use in pea-barley, faba bean-barley and lupin-barley intercrops grown at two temperate locations. *Journal of Agricultural Science* 142, 617–627. <https://doi.org/10.1017/S0021859604004745>.
- Krajewski P, D Chen, H Ćwiek, ADJ van Dijk, F Fiorani, P Kersey, C Klukas, M Lange, A Markiewicz, JP Nap, J van Oeveren, C Pommier, U Scholz, M van Schriek, B Usadel, and S Weise (2015). Towards recommendations for metadata and data handling in plant phenotyping. *Journal of Experimental Botany* 66, 5417–5427. <https://doi.org/10.1093/jxb/erv271>.
- Launay M, N Brisson, S Satger, H Hauggaard-Nielsen, G Corre-Hellou, E Kasynova, R Ruske, ES Jensen, and MJ Gooding (2009). Exploring options for managing strategies for pea-barley intercropping using a modeling approach. English. *European Journal of Agronomy* 31, 85–98. <https://doi.org/10.1016/j.eja.2009.04.002>.
- Lawler EL, JK Lenstra, and AHG Rinnooy Kan (1980). Generating All Maximal Independent Sets: NP-Hardness and Polynomial-Time Algorithms. *SIAM Journal on Computing* 9, 558–565. <https://doi.org/10.1137/0209042>.
- Licker R, M Johnston, JA Foley, C Barford, CJ Kucharik, C Monfreda, and N Ramankutty (2010). Mind the gap: How do climate and agricultural management explain the ‘yield gap’ of croplands around the world? *Global Ecology and Biogeography* 19, 769–782. <https://doi.org/10.1111/j.1466-8238.2010.00563.x>.
- Lobell DB, JM Deines, and SD Tommaso (2020). Changes in the drought sensitivity of US maize yields. *Nature Food* 1, 729–735. <https://doi.org/10.1038/s43016-020-00165-w>.
- Louarn G, R Barillot, D Combes, and A Escobar-Gutiérrez (2020). Towards intercrop ideotypes: Non-random trait assembly can promote overyielding and stability of species proportion in simulated legume-based mixtures. *Annals of Botany* 126, 671–685. <https://doi.org/10.1093/aob/mcaa014>.
- Louarn G, L Bedoussac, N Gaudio, EP Journet, D Moreau, E Steen Jensen, and E Justes (2021). Plant nitrogen nutrition status in intercrops– a review of concepts and methods. *European Journal of Agronomy* 124, 126229. <https://doi.org/10.1016/J.EJA.2021.126229>.
- Lowndes JSS, BD Best, C Scarborough, JC Afflerbach, MR Frazier, CC O’Hara, N Jiang, and BS Halpern (2017). Our path to better science in less time using open data science tools. *Nature Ecology & Evolution* 1. Number: 6 Publisher: Nature Publishing Group, 1–7. <https://doi.org/10.1038/s41559-017-0160>.
- Maat H (2011). The history and future of agricultural experiments. *NJAS - Wageningen Journal of Life Sciences* 57, 187–195. <https://doi.org/10.1016/j.njas.2010.11.001>.
- Mahmoud R, P Casadebaig, N Hilgert, L Alletto, GT Freschet, CD Mazancourt, and N Gaudio (2022). Species choice and N fertilization influence yield gains through complementarity and

- selection effects in cereal-legume intercrops. *Agronomy for sustainable development*. <https://doi.org/10.1007/s13593-022-00754-y>.
- Makowski D, T Nesme, F Papy, and T Doré (2014). Global agronomy, a new field of research. A review. *Agronomy for Sustainable Development* 34, 293–307. <https://doi.org/10.1007/s13593-013-0179-0>.
- Meunier C, L Alletto, L Bedoussac, JE Bergez, P Casadebaig, J Constantin, N Gaudio, R Mahmoud, JN Aubertot, F Celette, M Guinet, MH Jeuffroy, MH Robin, S Médiène, L Fontaine, B Nicolardot, E Pelzer, V Souchère, AS Voisin, B Rosiès, M Casagrande, and G Martin (2022). A modelling chain combining soft and hard models to assess a bundle of ecosystem services provided by a diversity of cereal-legume intercrops. *European Journal of Agronomy* 132. <https://doi.org/10.1016/j.eja.2021.126412>.
- Naudin C, G Corre-Hellou, S Pineau, Y Crozat, and MH Jeuffroy (2010). The effect of various dynamics of N availability on winter pea-wheat intercrops: Crop growth, N partitioning and symbiotic N-2 fixation. *Field Crops Research* 119, 2–11. <https://doi.org/10.1016/j.fcr.2010.06.002>.
- Naudin C, HMG van der Werf, MH Jeuffroy, and G Corre-Hellou (2014). Life cycle assessment applied to pea-wheat intercrops: A new method for handling the impacts of co-products. en. *Journal of Cleaner Production* 73, 80–87. <https://doi.org/10.1016/j.jclepro.2013.12.029>.
- Newman SJ and RT Furbank (2021). A multiple species, continent-wide, million-phenotype agronomic plant dataset. *Scientific Data* 8, 1–8. <https://doi.org/10.1038/s41597-021-00898-8>.
- Pelzer E, M Bazot, L Guichard, and MH Jeuffroy (2016). Crop Management Affects the Performance of a Winter Pea–Wheat Intercrop. en. *Agronomy Journal* 108, 1089–1100. <https://doi.org/10.2134/agronj2015.0440>.
- Phillips CA, K Wang, EJ Baker, JA Bubier, EJ Chesler, and MA Langston (2019). On Finding and enumerating maximal and maximum k-partite cliques in k-partite graphs. *Algorithms* 12. <https://doi.org/10.3390/a12010023>.
- Popkin G (2019). Data sharing and how it can benefit your scientific career. *Nature* 569, 445–447. <https://doi.org/10.1038/d41586-019-01506-x>.
- Senft M, U Stahl, and N Svoboda (2022). Research data management in agricultural sciences in Germany: We are not yet where we want to be. *PLOS ONE* 17. Ed. by Pulvento C, e0274677. <https://doi.org/10.1371/journal.pone.0274677>.
- Sparks AH (2018). nasapower: A NASA POWER Global Meteorology, Surface Solar Energy and Climatology Data Client for R. en. *Journal of Open Source Software* 3, 1035. <https://doi.org/10.21105/joss.01035>.
- Tang X, SA Placella, F Dayde, L Bernard, A Robin, EP Journet, E Justes, and P Hinsinger (2016). Phosphorus availability and microbial community in the rhizosphere of intercropped cereal and legume along a P-fertilizer gradient. *Plant and Soil* 407, 119–134. <https://doi.org/10.1007/s11104-016-2949-3>.
- Tannier, E. (2024). Collecting, assembling and sharing data in crop sciences. *Peer Community in Mathematical and Computational Biology*. 100197 <https://doi.org/10.24072/pci.mcb.100197>.
- Tardieu F (2020). Educated big data to study sensitivity to drought. *Nature Food* 1, 669–670. <https://doi.org/10.1038/s43016-020-00187-4>.
- Viguier L, L Bedoussac, EP Journet, and E Justes (2018). Yield gap analysis extended to marketable grain reveals the profitability of organic lentil-spring wheat intercrops. *Agronomy for Sustainable Development* 38, 39. <https://doi.org/10.1007/s13593-018-0515-5>.
- White JW and FK Van Evert (2008) Publishing agronomic data. *Agronomy Journal* 100,1396–1400. <https://doi.org/10.2134/agronj2008.0080F>.
- Wickham H (2014). Tidy Data. *Journal of Statistical Software* 59. <https://doi.org/10.18637/jss.v059.i10>.
- Wickham H and G Grolemund (2016). *R for data science: import, tidy, transform, visualize, and model data*. O'Reilly Media, Inc. <https://r4ds.had.co.nz/>

- Wilkinson MD, M Dumontier, IJ Aalbersberg, G Appleton, M Axton, A Baak, N Blomberg, JW Boiten, LB da Silva Santos, PE Bourne, J Bouwman, AJ Brookes, T Clark, M Crosas, I Dillo, O Dumon, S Edmunds, CT Evelo, R Finkers, A Gonzalez-Beltran, AJ Gray, P Groth, C Goble, JS Grethe, J Heringa, PA t Hoen, R Hooft, T Kuhn, R Kok, J Kok, SJ Lusher, ME Martone, A Mons, AL Packer, B Persson, P Rocca-Serra, M Roos, R van Schaik, SA Sansone, E Schultes, T Sengstag, T Slater, G Strawn, MA Swertz, M Thompson, J Van Der Lei, E Van Mulligen, J Velterop, A Waagmeester, P Wittenburg, K Wolstencroft, J Zhao, and B Mons (2016). Comment: The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data* 3, 1–9. <https://doi.org/10.1038/sdata.2016.18>.
- Wilson G, J Bryan, K Cranston, J Kitzes, L Nederbragt, and TK Teal (2017). Good enough practices in scientific computing. *PLOS Computational Biology* 13, e1005510. <https://doi.org/10.1371/JOURNAL.PCBI.1005510>.
- Zamir D (2013). Where Have All the Crop Phenotypes Gone? *PLoS Biology* 11, 1–4. <https://doi.org/10.1371/journal.pbio.1001595>.