

Insights into wheat science: A bibliometric review using unsupervised machine learning techniques

Martín Pérez-Pérez^{a,b,*}, Miguel Ribeiro^{c,d}, Florentino Fdez-Riverola^{a,b}, Gilberto Igrejas^{c,e,f}

^a CINBIO, Universidade de Vigo, Department of Computer Science, ESEI - Escuela Superior de Ingeniería Informática, 32004, Ourense, Spain

^b SING Research Group, Galicia Sur Health Research Institute (IIS Galicia Sur), SERGAS-UVIGO, Spain

^c Department of Genetics and Biotechnology, University of Trás-os-Montes and Alto Douro, 5001-801 Vila Real, Portugal

^d Chemistry Research Centre-Vila Real (CQ-VR), University of Trás-os-Montes and Alto Douro, 5000-801, Vila Real, Portugal

^e Functional Genomics and Proteomics Unit, University of Trás-os-Montes and Alto Douro, 5001-801 Vila Real, Portugal

^f LAQV-REQUIMTE, Faculty of Science and Technology, Nova University of Lisbon, 2829-516 Caparica, Lisbon, Portugal

ARTICLE INFO

Original content: [SING Group \(Original data\)](#)

Keywords:

Wheat
Literature analysis
Knowledge extraction
Machine learning
LLM
Clustering

ABSTRACT

Wheat (*Triticum* spp.) has been one of the most important cereal crops, serving as a source of protein and energy in the human diet. It remains a vital component of global food security, with extensive scientific literature dedicated to its study, although the large volume of literature often hinders global analysis. In this study, different unsupervised machine learning techniques, such as K-Nearest Neighbors (KNN) and Uniform Manifold Approximation and Projection (UMAP), text mining analyses, including word embeddings and statistical word analysis, and graph analysis methodologies, were applied to gain a deeper understanding of the wheat literature. The proposed bibliometric analysis was conducted and integrated with the Journal Citation Reports (JCR) to identify major wheat research trends in the PubMed literature. This analysis examined the evolution of these trends over time, evaluated the geographical distribution, impact, and research domains, and assessed author collaboration networks and the evolving relevance of different countries. Research on disease resistance, genetic modification, and dietary impact demonstrates a consistent increase in output, while interest in topics related to overcoming salt stress and enhancing animal feed appears to be diminishing. Interestingly, research on wheat germ agglutinin saw a surge in interest during the late 2000s, stabilizing thereafter. These trends underscore the dynamic nature of wheat research, driven by evolving priorities and technological advancements, particularly in genetics and omics tools. Moreover, the increasing significance of China in wheat research, including its size, impact, and networking, alongside longstanding leaders such as the United States, signals a shifting landscape in global wheat research.

1. Introduction

The cultivation of wheat has played a pivotal role in shaping human societies throughout history, leaving a significant mark on the biopsychosocial aspects of our existence. From its domestication in the Fertile Crescent to modern times, wheat has not only served as a staple food but has also shaped cultural practices, economic structures, and human nutrition and health (Ribeiro et al., 2013). Today, the most widely cultivated species is the allohexaploid wheat (*Triticum aestivum* L., $2n = 6x = 42$). In addition to common wheat, durum wheat (*Triticum durum* Desf.), an important tetraploid species used primarily for pasta

production (de Sousa et al., 2021). Spelt (*Triticum spelta* L.) and recognized for its nutritional value, has also gained increasing attention in recent years (Frakolaki et al., 2018).

In the modern era, the importance of wheat extends beyond sustenance, with the food industry being a major beneficiary. Wheat is a global staple and a primary ingredient in a diverse array of food products such as bread, pasta, noodles, and various other food items. When considering global production, certain countries stand out as major contributors to the cultivation of wheat. Wheat production, totalling approximately 808 million metric tons worldwide, is particularly abundant in countries like China (138 Mt), India (108 Mt), and Russia

* Corresponding author. CINBIO, Universidade de Vigo, Department of Computer Science, ESEI - Escuela Superior de Ingeniería Informática, 32004, Ourense, Spain.

E-mail addresses: martiperez@uvigo.gal, martiperez@uvigo.es (M. Pérez-Pérez), jmribeiro@utad.pt (M. Ribeiro), riverola@uvigo.gal, riverola@uvigo.es (F. Fdez-Riverola), gigrejas@utad.pt (G. Igrejas).

<https://doi.org/10.1016/j.jcs.2024.103960>

Received 4 April 2024; Received in revised form 14 June 2024; Accepted 17 June 2024

Available online 17 June 2024

0733-5210/© 2024 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

(104 Mt) (fao.org/faostat).

In wheat, gluten plays a crucial role in imparting distinct functional and technological characteristics to their flours. The term “gluten” is commonly used to refer to the complex protein mixture found in wheat, rye, or barley, being insoluble in water or salt solution. In wheat, gluten is made up of gliadins (prolamins) and glutenins (glutelins), while rye contains secalins, and barley contains hordeins (Wieser et al., 2023). Among these cereals, wheat stands out as the only known grain whose flours, when mixed with water and kneaded, can create a distinctive viscoelastic dough with high functional and technological value. This phenomenon is attributed to the formation of a complex three-dimensional protein network linked by disulfide bonds. Together, gliadins and glutenins, impart respectively viscosity and elasticity to wheat gluten (Ribeiro et al., 2019; Wieser et al., 2023).

However, it is the gluten that triggers severe immune-mediated diseases, such as celiac disease (Sabença et al., 2021). While its global prevalence remains uncertain, studies suggest a prevalence of 1%, therefore representing a public health problem. For example, in the United States, the prevalence of celiac disease has increased by a factor of 4–5 in the past 50 years. Specialists believe that this increase reflects a true rise in incidence rather than increased awareness and detection (Lebwohl et al., 2018; Ribeiro and Nunes, 2019). However, despite becoming one of the best-understood HLA-linked disorders, life-long adherence to a strict gluten-free diet is the only effective treatment available (Kupper, 2005).

Gluten-free diets, commonly prescribed for treating celiac disease and other gluten-related disorders, have become a growing trend as individuals without any gluten-related disorders are increasingly adopting gluten-free diets (Kim et al., 2016). As of February 19, 2024, a Google search for “gluten-free diet” yielded nearly 1 billion results. Many individuals opt for gluten-free foods under the belief that they are healthier than their gluten-containing counterparts (Sabença et al., 2021). However, there is no scientific evidence to support the notion that individuals without medical recommendations can benefit from this diet. In fact, it is sometimes considered an unhealthy option (Lebwohl et al., 2017).

In the context of the dynamic interplay between this cereal and humans, which translates into interdisciplinary research spanning agronomy, climate change, genetics, breeding, food science and technology, nutrition, and health, and the production of a vast literature corpus, the use of bibliometric analysis, facilitated by advanced language models, can emerge as an important tool.

In this sense, different unsupervised machine learning techniques, such as K-Nearest Neighbors (KNN) and Uniform Manifold Approximation and Projection (UMAP), along with text mining techniques, containing word embeddings and statistical word analysis, as well as graph analysis methodologies were applied to gain a deeper understanding of the wheat literature. The proposed bibliometric analysis was employed and integrated with the journal citation report to identify major wheat research trends in the PubMed literature. The purpose of the present work is to perform a quantitative bibliometric analysis of the domain literature within the “wheat” domain from 1980 to 2023. This includes evaluating countries, journal metrics, and research fields to assess research impact and trends, as well as identifying patterns in author collaborations and emerging topics related to health and agriculture, given the critical importance of wheat in both human nutrition and agricultural economies. Understanding the research landscape allows for pinpointing advancements in areas like the nutritional enhancements of wheat, its role in preventing health disorders, and innovations in sustainable agricultural practices.

2. Materials and methods

This section describes the proposed statistical and unsupervised data analysis workflow applied to the cereals-related Bibliome to extract and evaluate the literature knowledge.

The term “Bibliome” refers to the comprehensive and structured collection of bibliographic data and literature content in a specific biomedical domain (Alfred, 2001; Searls, 2001), in this case, wheat. Thus, Fig. 1 depicts the pipeline followed in the current study to obtain, process, and analyse the cereal-related scientific Bibliome, while the following subsections detail the different decisions applied.

In this study, wheat-related abstracts and author metadata were extracted from PubMed to conduct a comprehensive bibliometric analysis. The extracted data underwent several preprocessing steps: authors were reconstructed as collaboration graphs, author metadata were processed to uncover demographics and institutional information, and article abstracts were transformed into embeddings and clustered to identify the main topics discussed, using statistical word analysis methods such as term frequency-inverse document frequency (TF-IDF). This approach facilitated the extraction and reconstruction of valuable bibliometric knowledge from the literature. The following sections describe this process in greater detail.

2.1. Data retrieval

As illustrated in Fig. 1, the objective of this phase was to retrieve domain-related scientific literature and their metadata from the PubMed repository, and to vectorize the text content of the different articles to support the following phase.

2.1.1. Bibliome and metadata assembly

To carry out the proposed knowledge extraction analysis, a collection of 114,748 PubMed abstracts that mentioned the word “Wheat” in their title or abstract from 1980 to 2023 were collected using The Entrez Utilities Web services from the NCBI (National Center for Biotechnology Information). Along with the abstract and their title, metadata such as authors and their affiliations, funding information, keywords, MESH terms, journal names, and publication dates were also collected.

Moreover, the following phases integrated complementary knowledge databases such as (i) the Journal Citation Report (JCR) between 2008 and 2022, which contains the journal impact factor, categories and quartiles for evaluating and contextualizing the relevance of publications; and (ii) the GeoNames database (Geonames.org, 2016), which contains more than 10 million topographical names. JCR is a widely recognized resource for evaluating and comparing academic journals (Salisbury, 2020). Published by Clarivate Analytics, JCR provides a systematic, objective means to critically assess the world’s leading journals, with quantifiable, statistical information based on citation data. Key metrics provided by JCR include the Quartile classification, which reflects a ranking system used within JCR to categorize journals into four quartiles (Q1, Q2, Q3, and Q4) based on their impact factor within a specific subject category.

2.1.2. Bibliome vectorization

To vectorize the text content of the different articles, a pre-trained large language model (LLM) was used. Recently, models such as GPT-3 and BERT have significantly improved text analysis by generating intricate text vectors, enabling robust comparison and sophisticated clustering of documents based on semantic similarity within a vector space. For this reason, the pre-trained scientific-domain specific BERT-based model (PubMedBert) proposed by Gu Y et al. (Gu et al., 2022) was applied.

This selection offers distinct advantages within the researched domain over general language models such as GPT or standard BERT. PubMedBert, being specifically trained in a vast corpus of biomedical literature, is finely tuned to the intricacies and specialized terminology found in this field. This specialization allows for a higher degree of accuracy and relevance in semantic analysis and clustering tasks within biomedical and related scientific texts. Additionally, as a model in the BERT family, its bidirectional context-aware embeddings enable it to capture rich contextual information from both the left and right contexts

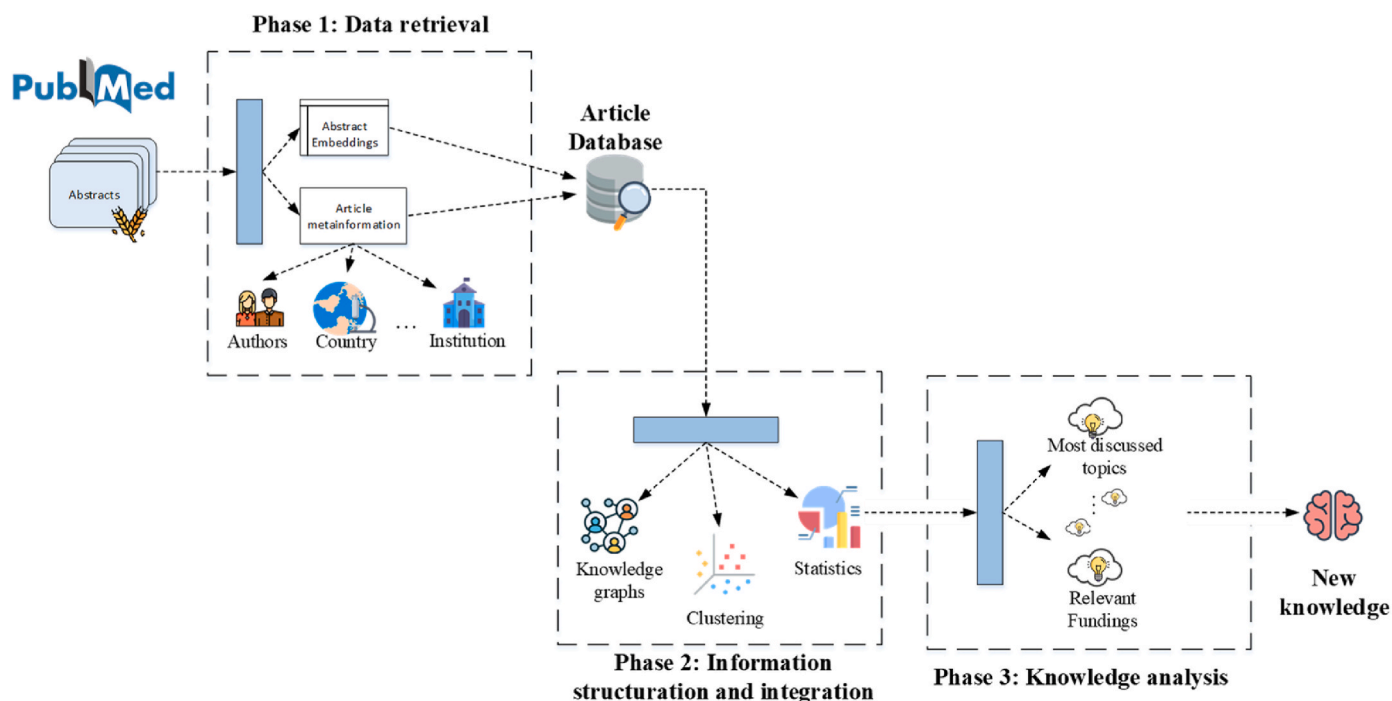


Fig. 1. Workflow for Data Extraction, Analysis and Visualization of the domain scientific Literature.

of a word. In contrast, other model architectures, such as GPT models make it less preferable for classification or clustering tasks due to their unidirectional, processing of text from left to right.

Therefore, PubMedBert embeddings were employed to extract text features and group similar articles together. At the end of this process, articles (titles and abstracts) were transformed into a numerical vector representation, often referred to as document embeddings.

2.2. Information structuration and integration

This phase aimed to use various machine learning, statistical, and NLP methods to organize and standardize article information, extracting new insights from the bibliome, with subsequent sections detailing the implemented approaches.

2.2.1. Topic extraction

In order to extract the most discussed topics in the evaluated domain, a combination of dimensionality reduction, unsupervised clustering techniques and concept statistical methods were applied.

2.2.1.1. Bibliome dimensionality reduction. To obtain the different topics discussed within the evaluated domain, it was necessary to apply different clustering techniques to discover the nearest articles, that is, discover the document embeddings with the nearest vector space. However, given that working with such high-dimensional embeddings may be computationally and memory-intensive, a dimension reduction strategy was applied. To the best of the authors' knowledge, the UMAP (McInnes et al., 2018) technique has the ability to capture non-linear patterns and preserve local structures, thus making it a suitable choice for reducing LLM embedding dimension. UMAP generates a lower-dimensional projection that preserves the similarities and dissimilarities between data points in high-dimensional space by taking their closeness into account (Leetaru, 2024).

2.2.1.2. Clustering: group the main discussion topics. With the aim of grouping together articles based on their content and extract the main topics of discussion in the current domain, the word vector embedding representation of the articles enables the application of clustering

techniques such as KNN (Li et al., 2021). To carry out this process, the KNN algorithm calculates the distance or similarity between document embeddings to determine the nearest neighbors based on the similarity of its LLM vector representation to the centroid or representative points of the cluster. To select the optimal number of clusters, the popular Elbow Method was applied in combination with the domain knowledge and expertise analysis of the final topics (Makwana et al., 2013). At the end of this process, document embeddings were grouped into different clusters representing the different discussed topics.

2.2.1.3. Word relevance and clustering representation. Once the documents were grouped based on their vector representation, the next step was to identify the most commonly discussed topics using a word statistical method. TF-IDF (Term Frequency-Inverse Document Frequency) is a common representation method used in text analysis to represent the importance of a word in a document concerning a corpus of documents. The TF measures the frequency of a word in a document, while the IDF measures the rarity of a word in the dataset. In the context of clustering, the TF-IDF representation was used to identify the keywords or phrases that are characteristic of a particular cluster and provide a more interpretable and concise representation of the domain key topics in discussion. In this way, TF-IDF is defined as follows (Equation (1)):

$$tf - idf(t, s, D) = tf(t, s) \times idf(t, D) \quad (\text{Eq 1})$$

where t is the evaluated term, s stands for any given sentence of the dataset, D , and $tf(t, s)$ expresses the ratio corresponding to the term, t , in a sentence, s , described as follows (Equation (2)):

$$tf(t, s) = \frac{n_t}{\sum_k n_k} \quad (\text{Eq 2})$$

where n_t is the number of occurrences of the term, t , in a sentence, s , and n_k is the total number of terms in a sentence, s . Moreover, in Equation (1) $idf(t, s)$ stands for the logarithmic ratio of the term, t , in the dataset, D , and is computed as follows (Equation (3)):

$$idf(t, D) = \log \frac{|D|}{|\{s_i \in D \mid t \in s_i\}|} \quad (\text{Eq 3})$$

At the end of this process, document clusters were represented by the most relevant discussed concepts (from 1 to 3 grams).

2.2.2. Knowledge normalization

2.2.2.1. Semantic normalization. The process of semantic normalization involved different steps: (i) word tokenization and n-gram generation to segment the text into semantically significant units like words and phrases; (ii) frequent and low-utility terms removal, including common English stop words, words with fewer than two characters or containing numerical digits and domain-specific keywords that may be considered as stop words; and finally in the last step (iii) tokens were normalized applying lemmatization techniques to represent concepts with the lexeme form. Moreover, the sequence of n-grams is considered in identifying tokens as a unified concept; for instance, “wheat flour” and “flour wheat” are regarded as identical terms.

2.2.2.2. Author affiliation normalization. Author affiliations are often written manually and freely so obtaining the affiliation country of authors can be a difficult and tedious task. Therefore, to identify, extract and normalize the author affiliation country, a three-step pipeline was applied in cascade and if no countries or different countries were identified, the location was set to “Unknown” and the subsequent step was applied. The steps were applied in the following order: (i) all author affiliations were parsed using the Country named entity recognition Python library (Wood, 2022); (ii) all remaining “Unknown” affiliation locations were searched at the GeoNames database; (iii) all remaining locations were prompted in the Opensource chatbot Vicuna to extract the country (Chiang et al., 2023) and finally; (iv) all remained “Unknown” locations were searched against the Geopy limited API (“Geopy,” 2023). Consequently, 128,207 affiliation countries of a total of 128,567 affiliations (i.e., 96.6%) were identified with an F. score of 0.99 in a random sample of 1000 entities. Many of the recognition errors were induced by the PubMed retrieval of multiple affiliations within the same record.

2.2.3. Knowledge integration: research areas and impact factor

In order to evaluate the different research fields in the domain and the relevance of the different studies, the JCR statistics of the year 2022 were integrated considering the ISBN/ISS of the journal of publication. Articles from 1980 to 2023 were measured based on this index report (Salisbury, 2020).

2.2.4. Knowledge analysis metric establishment

The relevance of entity publications (e.g., countries or affiliation entities) can be measured from two points of view, (i) as cumulative numbers, to measure the entities that have more relevance in terms of production in the studied domain, or (ii) the average relevance or impact factor, which serves to measure the effectiveness of the studies. In this sense, this metric can also evaluate the research impact of countries with smaller populations, which inherently have fewer opportunities for publication and, consequently less Bibliome production.

Therefore, the following equations (Equations (4)–(6)) formally define the rationale behind the evaluation of the literature from both points of view.

$$\text{Quartile ratio} = \frac{\sum a_{Qx}}{\sum a} \quad (\text{Eq 4})$$

Where a represents an article published in the studied domain and a_{Qx} represent an article in the Qx quartile (i.e Q1, Q2, Q3 or Q4)

$$\text{Average impact factor} = \frac{\sum a_{IF}}{\sum a} \quad (\text{Eq 5})$$

Where a represents an article published in the studied domain and a_{IF} represent the JCR article impact factor

$$\text{Publication percentage} = \frac{\sum a_k}{\sum a} \quad (\text{Eq 6})$$

Where a represents an article published in the studied domain and a_k an article in a specific k cluster (main discussed area).

2.2.5. Formal collaboration knowledge structuration

In order to define what is considered a country collaboration, this section formally defines an article relationship of two countries to generate the interrelation country knowledge graph. In this sense, a vertex v in the country collaboration graph is defined as the following (Equation (7))

$$v = \{E_1, E_2 \dots E_n\} \quad (\text{Eq 7})$$

Where E_{th} represent the affiliation entities of the same country and is defined as: (Equation (8))

$$E = \{a_1, a_2 \dots a_m\} \quad (\text{Eq 8})$$

Where a_{th} represent the i -th article published in the studied domain by the entity or affiliation.

In terms of country interrelations, an edge e in the country collaboration graph is defined as the following (Equation (9)):

$$e_{v_1v_2} = \{a_{v_1v_1}, a_{v_1v_2} \dots a_{v_1v_3}\} \quad (\text{Eq 9})$$

Where $a_{v_1v_{th}}$ represent the i -th article and their metadata in which the vertex v_1 and v_2 participates. This graph formalization methodology enables the application of visualization and analysis techniques such as the edge and vertex degree proportionality based on: (i) the number of articles that conform the graph or (ii) the impact factor of the publications; as well as (iii) other state-of-the-art graph metrics such as the betweenness centrality, degree, closeness centrality or clustering coefficient, among others (Bolland, 1988; Nieminen, 1974).

3. Results & discussion

3.1. Publications and trends over time

In order to revise the literature volume evolution, the next figure represents the volume of publications from 1980 to 2023 related to the topic “wheat” (Fig. 2A). Notably, it reveals that the number of publications per year approximately follows an exponential line and, with a substantial surge in the number of publications, particularly evident after the year 2000.

Giraldo et al. (2019), found similar results in a bibliometric comparative analysis between wheat and barley, in particular since the second half of the twentieth century. This exponential growth is making it increasingly necessary to apply unsupervised techniques and data mining to obtain and structure latent knowledge in publications. This underscores the significance of works such as the current one since the information overload in the Bibliome already exceeds the ability of researchers to digest it and is growing at an unprecedented rate (Lyson et al., 2019).

In order to better understand the research that supports these results we have conducted a further analysis. Fig. 2B represents the most relevant terms in the topic “wheat” considering their TF-IDF (Equation (1)) in each identified main discussion area (i.e. article cluster) following the proposed analysis workflow. To obtain these results, only abstracts were analysed. Therefore, for the classification and naming within the clusters—each representing a collection of papers discussing related topics—Fig. 2B displayed the terms that occurred most frequently in each group. The relevance of these terms diminished progressively from left to right. For example, in a specific cluster, the primary terms identified were “Gene expression,” “Wheat Straw,” and “Winter wheat.”

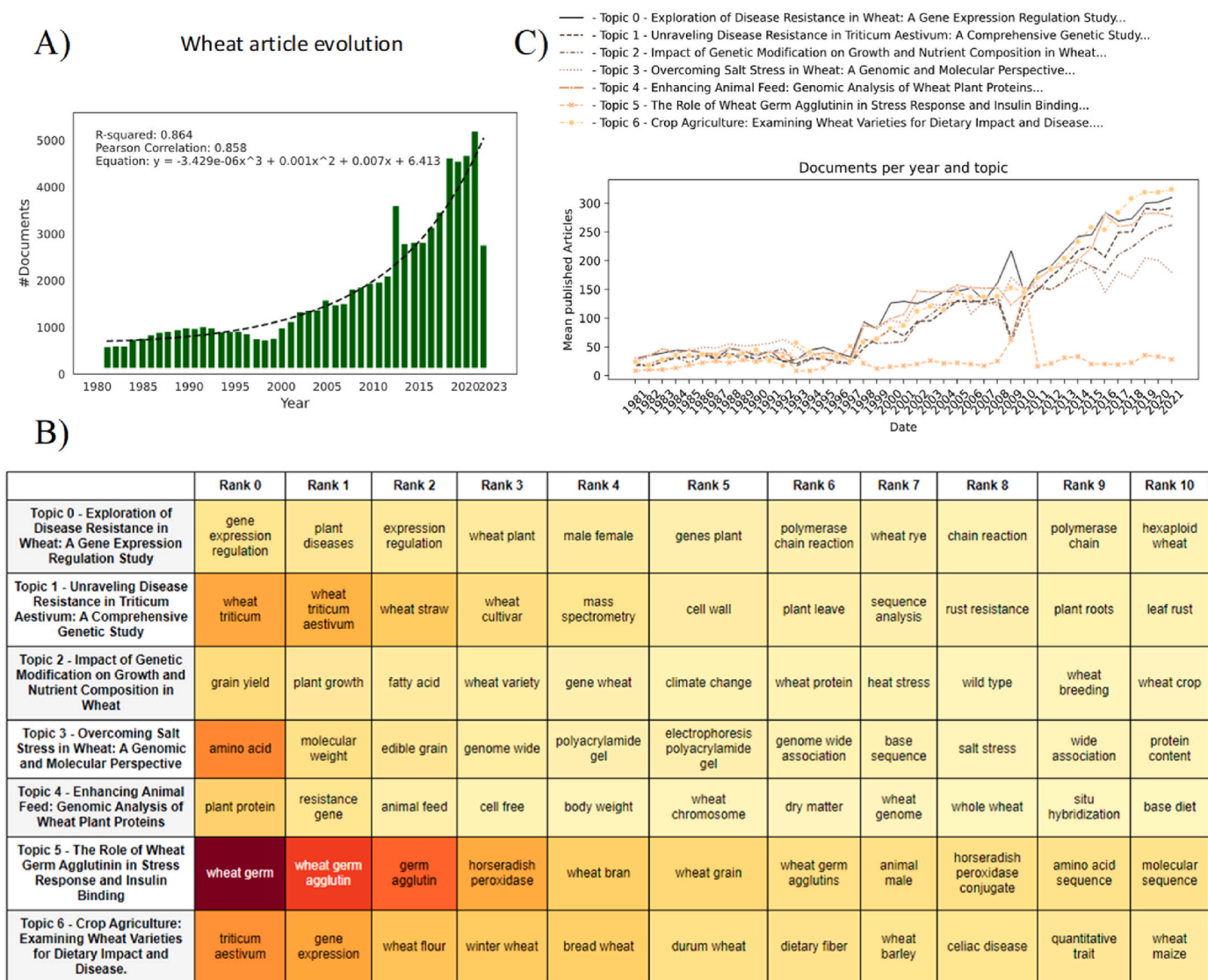


Fig. 2. Literature volume evolution. Number of publications per year from 1980 to 2023 related to the topic “wheat” (A). **The literature discussed topics (B).** The first column represents the name of the cluster or topic based on the TF-IDF of the discovered concepts. Columns from two to last, represent the relevance of the word ranked by their TF-IDF (i.e., Rank 0 is the most relevant word of the cluster). Background heatmap colour represents the relevance of the word inside the cluster. **Publication evolution categorized by identified clusters over the years (C).** Depicts the publication volume of each identified cluster from 1980 to 2023. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Based on these key terms, titles for each cluster were manually crafted, ensuring that the most significant terms had a prominent influence on the selection of each title. The effectiveness of the derived titles depends on various parameters and word evaluation methods, such as TF-IDF and simple occurrence counts. Each topic contains the most representative words and the same word cannot be in more than one topic; instead, it is selected where it is most relevant.

Although the trends identified in the research topics reflect diverse areas of study discussed from 1980 to 2023 in wheat research (Fig. 2B), it should be noted that genetics/genomics appears in 5 of the 7 identified topics. This observation elucidates the significant focus on genetic mechanisms and molecular-level understanding within wheat research endeavours, including advancements in wheat breeding, disease resistance studies, and biotechnological interventions aimed at enhancing crop productivity and resilience (Ma, 2023).

Considering both the topics and terms, the topic “Exploring disease resistance in wheat: a gene expression regulation study” underscores the importance of understanding the molecular mechanisms underlying disease resistance in wheat. Gene expression regulation ranked as the

top term, highlighting the focus on identifying specific genes and pathways involved in the plant’s response to diseases, which is crucial for developing resistant wheat varieties to combat plant diseases, as indicated by the term “Plant diseases” (Jabran et al., 2023; Jost et al., 2023). The second topic continues the emphasis on disease resistance but expands to a broader genetic investigation of wheat, particularly focusing on *Triticum aestivum*, which is the most important wheat species in terms of the economic importance on the consumer side. It emphasizes genetic traits specific to wheat, examines broader aspects of wheat physiology, and incorporates advanced analytical techniques (e.g. mass spectrometry). In fact, mass spectrometry is being used as a fundamental technique in proteomics and helps to fill the gaps in knowledge related to the dynamic interplay between genome, transcriptome and proteome (Ribeiro et al., 2013). Also, the topic delves into genetic diversity among wheat cultivars and investigates resistance mechanisms against rust pathogens, potentially incorporating physiological analyses of committed tissues (e.g. plant leaves). Rust pathogens pose significant threats to wheat crops worldwide, causing devastating yield losses and economic damage. With their ability to rapidly spread and infect wheat

plants, rust diseases, including leaf rust, stem rust, and stripe rust, are among the most concerning challenges in wheat production. Effective management strategies, including breeding for rust resistance, are crucial for safeguarding wheat yields and ensuring global food security (Jost et al., 2023).

The topic “Impact of genetic modification on growth and nutrient composition in wheat” delves into the realm of biotechnology and its implications for wheat breeding. The prominence of terms like “Grain yield” and “Plant growth” reflects the primary focus on enhancing agronomic traits to increase yield and improve plant vigour through genetic modification. This is one of the main focuses of research in wheat since increasing yield suggests a higher harvest value concerning the allocated area. This research focus aligns with the overarching goal of wheat research, which centres on enhancing crop productivity to meet the demands of a growing global population (Araus et al., 2019). By improving agronomic traits such as grain yield and plant growth through genetic modification, researchers aim to achieve higher harvest values relative to the allocated agricultural area, and concerns about climate change (Gupta, 2024). Increasing wheat yield is crucial for ensuring food security and economic stability, particularly in regions where wheat serves as a staple food crop. Higher yields not only contribute to meeting the dietary needs of growing populations but also enhance the profitability and sustainability of agricultural systems. Additionally, terms such as “Fatty acid” indicate a specific interest in modifying nutrient composition for improved nutritional quality and health benefits. For example, increasing the levels of omega-3 fatty acids can enhance the cardiovascular health benefits of wheat-based products. Conversely, reducing the levels of unhealthy saturated fats or trans fats can improve the overall healthfulness of wheat-based foods (Amjad Khan et al., 2017).

“Overcoming salt stress in wheat: a genomic and molecular perspective” addresses a critical challenge in wheat production—salt stress. Salt stress poses a severe threat to wheat productivity, particularly in regions with high soil salinity. The inclusion of terms like “Amino acid” and “Molecular weight” suggests a focus on understanding the biochemical and molecular mechanisms involved in salt tolerance. This research aims to identify genes and pathways associated with salt tolerance to develop salt-resistant wheat varieties, as seems to be suggested by the term “Edible grain” (Amirbakhtiar et al., 2021).

The topic “Enhancing animal feed: genomic analysis of wheat plant proteins” indicates applications beyond direct human consumption, emphasizing the importance of wheat in animal feed production. The term “Plant protein” highlights efforts to enhance protein content in wheat for animal feed applications, while terms like “Resistance gene” and “Animal feed” indicate potential genetic strategies to improve feed quality and animal health (Yang and Shen, 2018).

There is considerable literature on “The role of wheat germ agglutinin in stress response and insulin binding” which reflects the topic’s significant scientific interest and relevance. Nevertheless, the number of publications is on average significantly lower over the years in relation to the other topics identified. Wheat germ agglutinin is a well-characterized wheat protein renowned for its affinity for N-acetylglucosamine. This specific binding capacity gives wheat germ agglutinin diverse properties, including antifungal and cytotoxic effects, as well as demonstrated anticancer properties against various cancer cell types (Balčiūnaitė-Murzienė and Dzikašas, 2021). Additionally, the role of wheat germ agglutinin in enhancing insulin binding is also mentioned. It can accelerate the rate of insulin-receptor complex formation, without affecting the rate of dissociation of the insulin-membrane complex or the total number of insulin-binding sites (Cuatrecasas, 1973). The involvement of wheat germ agglutinin in plant stress response mechanisms suggests a focus on understanding its role in enhancing plant resilience and disease resistance under adverse conditions, highlighting its significant scientific interest and relevance.

Lastly, “Crop agriculture: examining wheat varieties for dietary impact and disease” underscores the interdisciplinary nature of wheat

research, integrating agronomic, nutritional, and health-related perspectives. Terms such as “*Triticum aestivum*,” “Wheat flour,” and “Dietary fiber” indicate a holistic examination of wheat varieties’ dietary attributes, including nutritional composition and potential implications for human health. Additionally, terms like “Gene expression” and “Celiac disease” suggest a comprehensive investigation into genetic factors influencing dietary impact and disease susceptibility, aiming to develop wheat varieties that meet both nutritional and health-related needs (Ribeiro and Nunes, 2019; Shewry and Ward, 2012). Intraspecific variation in wheat is researched to find varieties with desired traits, notably those with varying levels of celiac disease-related epitopes, indicating the possibility of selecting varieties with reduced toxicity for celiac disease sufferers. (Ribeiro and Nunes, 2019).

In order to evaluate the time evolution of the different literature-discussed topics, Fig. 2C depicts the publication evolution categorized by the identified clusters over the years. For topics 0, 1, 2, and 6, which focus on exploring disease resistance, genetic studies, genetic modification, and crop agriculture, respectively, there has been a notable increase in research output over time. Of special note is the key role of genetics, genomics, gene expression, and other omics tools, which are transversal to the first three topics, probably reflecting the continuous efforts to address critical challenges in wheat cultivation, such as combating diseases, improving genetic traits, enhancing nutritional composition, and evaluating varieties for agricultural suitability (Ma, 2023; Shewry and Ward, 2012). Topic 6, which boasts a higher number of papers, reflects an alignment with changing consumer preferences and health recommendations. This trend emphasizes the growing consideration for healthier dietary alternatives and increased awareness of diseases (Sabença et al., 2021).

On the other hand, topics 3 and 4, which pertain to enhancing animal feed and overcoming salt stress in wheat, respectively, appear to be experiencing a decline in research output. This decrease could be attributed to several factors, including saturation of research in these areas, technological advancements that have mitigated some of the challenges associated with animal feed and salt stress, or shifts in research priorities towards other pressing agricultural issues.

Research on Topic 5, which examines the role of wheat germ agglutinin in stress response and insulin binding, significantly increased from 2008 to 2010, then stabilized. The initial surge was likely due to biotechnological advancements, enabling deeper study into its properties and biomedical applications, as indicated by (Balčiūnaitė-Murzienė and Dzikašas, 2021). The continued stable research output reflects sustained interest in wheat germ agglutinin’s varied roles in plant physiology and biomedical sciences.

3.2. Origin, impact and type of research

Considering the origin, authors, and impact of the data, we conducted additional analyses based on the country of origin and the quartile of publications. The next figure (Fig. 3A) represents the total publications per quartile in the domain by the affiliations of the authors that are associated with the publication. Regarding the countries with the highest number of publications within the domain, China, the United States of America, Australia, and Germany, the first two among the three highest population rates in the world (United Nations- DESA, 2022; Worldometer, 2022), had the highest number of publications (i.e. production) in the 1st quartile. It should be noted that India, the second country with the biggest population density (United Nations- DESA, 2022; Worldometer, 2022), is one of the countries that published the most in the Q2 and Q3 quartiles, which has not yet achieved the same success rates in the Q1 quartile.

China, the world’s largest wheat producer (fao.org/faostat), has emerged as a major player in wheat research due to its substantial investment in agricultural science and technology. The country’s large population, coupled with the importance of wheat as a staple food crop, drives extensive research efforts aimed at improving wheat varieties,

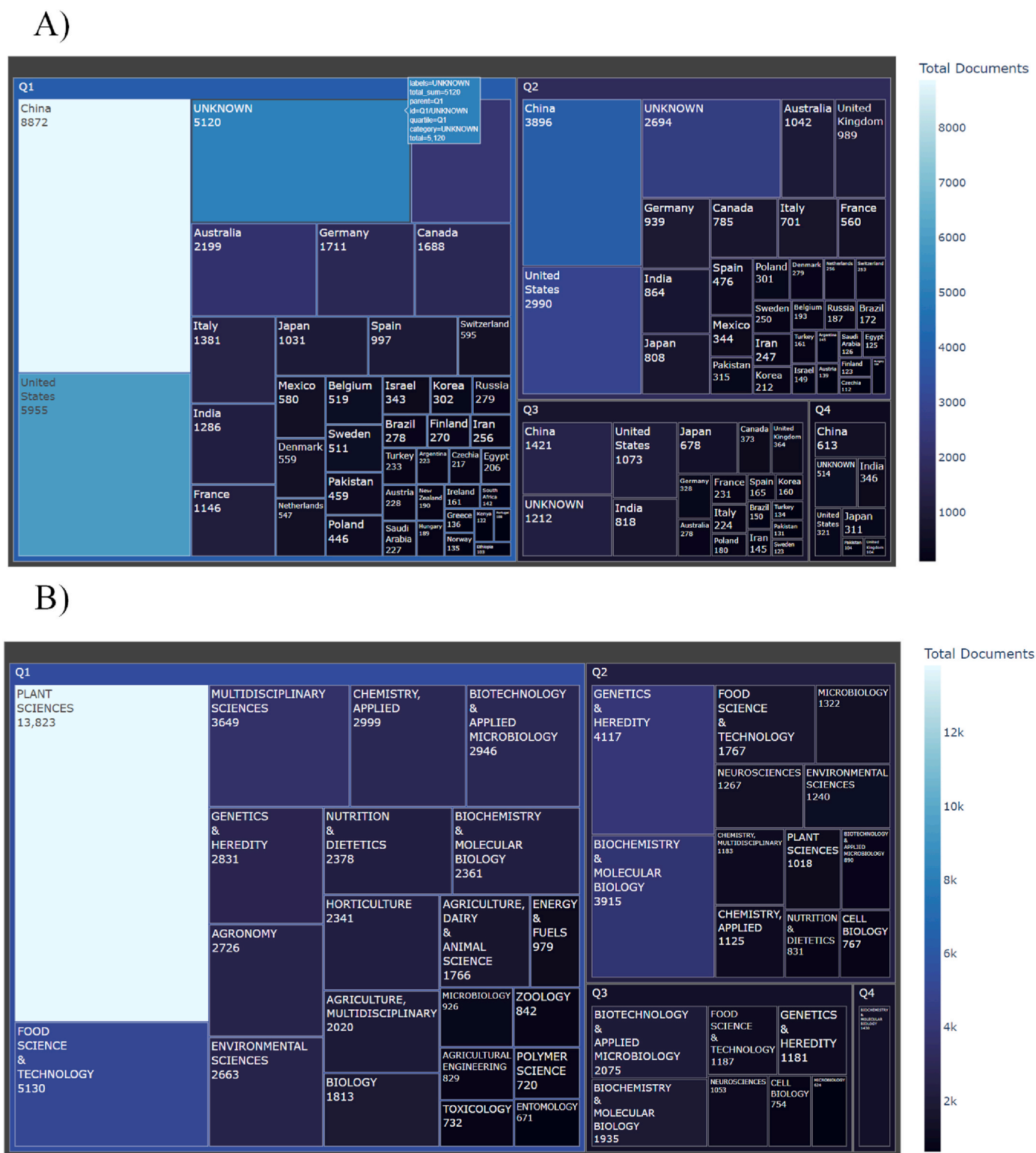


Fig. 3. Country treemap of the topic “Wheat” from 1980 to 2023 grouped by quartile and ranked by the number of total publications (A). Depicts the relevance of each country in each JCR quartile taking into account the total number of publications in the evaluated domain with more than 200 publications. The country was extracted from the affiliations of the authors that are associated with the publication. The most popular knowledge category groups of the topic “Wheat” from 1980 to 2023 were grouped by quartile and ranked by total publications (B). Depicts the relevance of each JCR category group (e.g. “Chemistry applied”) in each JCR quartile with more than 600 publications. That is the most popular JCR category group by quartile.

enhancing yield and quality, and addressing challenges such as disease resistance (Qin et al., 2015).

Considering the total number of publications by quartile and JCR category (Fig. 3B), the areas of Plant Sciences, with a much higher

number than all others, and Food Science & Technology stand out as the most important categories. The importance of wheat lies substantially in its importance as food for humanity (Shewry, 2009). In this sense, studies seek to “improve” the plant and the entire spectrum related to it

as food. In fact, the topics identified previously align with these categories both in the search for disease resistance, genetic traits of agronomic interest, and those associated with growth and yield, as well as in the impact on the human diet and diseases. Furthermore, this reflects agricultural advancement and its harmonization with nutritional and environmental considerations in terms of research, which is in line with (Giraldo et al., 2019). Overall, the convergence of “Multidisciplinary Sciences”, “Chemistry, applied,” “Biotechnology & Applied Microbiology” and other transversal areas, highlights the collaborative and interdisciplinary nature of research in unlocking the potential of cereals for addressing global challenges.

3.3. Research network

The next figure represents the top collaborations of the publication countries taking into account the affiliation of the authors (Fig. 4).

The analysis of the country collaboration graph delineates a decentralized configuration, exemplified by an average collaboration index of 13.48 among nations, highlighting extensive global research partnerships. Predominantly, the United States, China, and Australia are identified as central nodes, reflecting their significant roles due to their robust research infrastructure and prolific output. The network

architecture, characterized by a radius of two and a diameter of four, suggests a balanced distribution without overarching dominance by any single entity, as denoted by a centralization score of 0.550. Furthermore, the network’s structure reveals pronounced clustering, evidenced by a high clustering coefficient of 0.833 and a moderate density of 0.188, indicating the presence of cohesive subgroups, potentially organized around geographic proximities or shared academic fields.

The analysis further identifies China, the United States, Germany, Australia, the United Kingdom, Canada, and India as prominent nations within the network, demonstrating extensive direct collaborations and central roles in the global research nexus. This complex network of relationships is crucial for facilitating comprehensive cross-border collaborations and knowledge dissemination. Notably, countries such as Russia, Mexico, and Spain, the first two among the top ten highest population rates in the world (United Nations- DESA, 2022; Worldometer, 2022), are highlighted for their significant betweenness centrality, acting as vital intermediaries in the research community and enhancing network cohesion.

This intricate pattern of global research collaboration underscores the dynamic nature of international academic interactions, characterized by both broad-based and specialized partnerships. Detailed exploration and additional insights into this analysis are available in

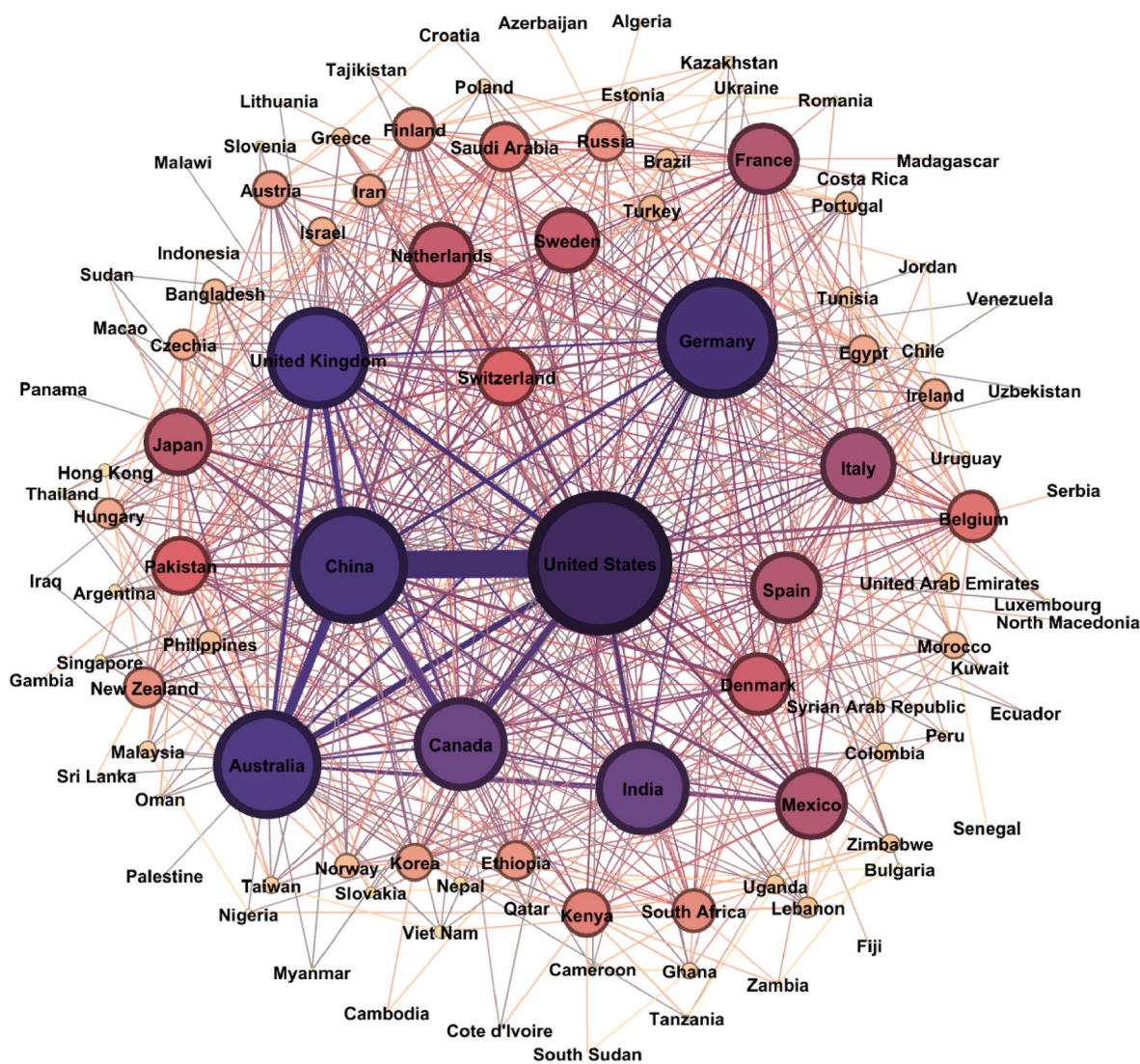


Fig. 4. Top article collaborations map. Graph node size and colour depict the relevance of the country taking into account the total number of collaborations with other countries, whereas the edge colour and weight represent the total number of collaborations between both countries. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Supplementary Material 1, providing an expanded perspective on the evolving landscape of global scientific collaboration.

3.4. Country relevance evolution

Analysing the last 43 years of publication volume, China and India have emerged as the most prominent nations since the turn of the new millennium. The case of China is particularly noteworthy since it has, within a period of only 10 years, reached a volume of publications equal to total number produced in 2011 by the United States of America, the main protagonist in the previous 31 evaluated years. After this, the volume of publications accumulated in the United States is beginning to lose prominence in favour of China (as indicated in the chart). However, it should be noted that China and India are the two countries with the highest population density in the world today (United Nations- DESA, 2022; Worldometer, 2022).

Conversely, when considering not only the publication volume but also the average impact factor of these publications, Fig. 5 describes the accumulated relevance of the countries across distinct decades. In this case, only countries with more than 50 publications per year are

considered. This perspective enables the discretion of when and in what dimension countries were integrating themselves into the research landscape of this domain.

Fig. 5 depicts how the contributions of different countries to research have evolved over time, showing an initial dominance by the United States, China, Australia, Canada, and the United Kingdom until 2013. From 2013 onwards, India and various European countries started to emerge as significant contributors. By the end of 2018, China's academic influence increased, indicating a shift towards greater global participation in the field, while the influence of the United States declined. By 2022, China further solidified its leading position, with Canada and Spain also becoming more prominent, particularly in producing top researchers. This suggests an ongoing trend towards globalization in academic research. Despite varying publication volumes, the United States, China, Australia, Canada, and India remain the primary sources of significant research in recent years, as highlighted by both Figs. 5 and 4.

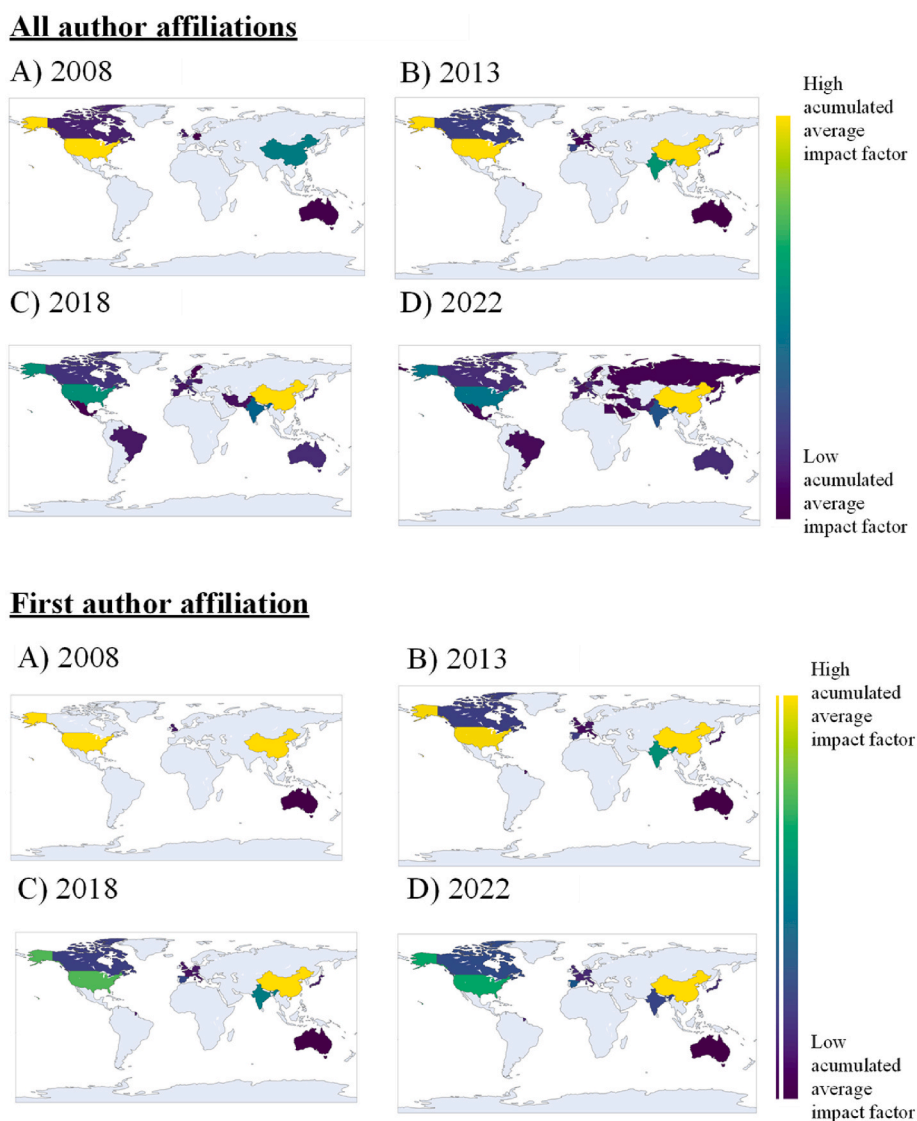


Fig. 5. Country publication's evolution from 2008 to 2022 ranked by their average impact factor accumulated. Only countries with more than 50 annual publications were evaluated. Countries were slowing down as the average impact factor cannot be maintained taking into account the total volume of publications. A) Represents the accumulated average impact factor from 2008; B) Represents the accumulated average impact factor from 2013; C) Represents the accumulated average impact factor from 2018; and D) Represents the accumulated average impact factor from 2008 to 2022.

4. Conclusions

Wheat was a very important food for primitive societies and continues to be one of the most important foods for humans today, which is reflected in the size of scientific literature targeting wheat as a main topic of interest. The major trends in wheat research identified encompass a diverse array of topics, each contributing valuable insights to the field. Those focusing on disease resistance, genetic modification, and dietary impact, have exhibited a growing trend in research output over time. This trend underscores the ongoing efforts to address critical challenges in wheat cultivation, such as combating diseases, improving genetic traits, and enhancing nutritional quality. Nevertheless, topics related to overcoming salt stress and enhancing animal feed, appear to be losing traction in terms of research interest. Wheat germ agglutinin-related research experienced a significant surge in research activity between 2008 and 2010, which is probably related to its interest in biomedical applications but has since remained relatively stable. These trends reflect the dynamic nature of wheat research, driven by evolving priorities and technological advancements. Moreover, the noteworthy role of research grounded in genetics or omics tools within the primary research trends is striking. The rapid advancement of genomic technologies since the early 2000s has revolutionized cereal research. The completion of genome sequencing projects for major cereal crops such as rice, maize, wheat, and barley has provided researchers with a wealth of genetic information, enabling them to identify key genes associated with important traits such as yield, disease resistance, and nutritional quality. On the other hand, the increasingly important role that China is playing in the field of wheat research must be highlighted. Once marked by other important players (with continued prominence), such as the United States, China now has a very relevant role both in the number and impact of research and in the number of global collaborations.

CRediT authorship contribution statement

Martín Pérez-Pérez: Writing – review & editing, Writing – original draft, Data curation, Conceptualization. **Miguel Ribeiro:** Writing – review & editing, Writing – original draft, Methodology, Conceptualization. **Florentino Fdez-Riverola:** Writing – review & editing, Investigation, Conceptualization. **Gilberto Igrejas:** Writing – review & editing, Validation, Conceptualization.

Declaration of competing interest

The authors declare that they have no conflicts of interest.

Data availability

No data was used for the research described in the article.
[SING Group \(Original data\)](#) (SING Group)

Acknowledgements

SING group thanks CITI (Centro de Investigación, Transferencia e Innovación) from the University of Vigo for hosting its IT infrastructure. This work was supported by: the Associate Laboratory for Green Chemistry - LAQV financed by the Portuguese Foundation for Science and Technology (FCT/MCTES) [LA/P/0008/2020 DOI 10.54499/LA/P/0008/2020, UIDP/50006/2020 DOI 10.54499/UIDP/50006/2020 and UIDB/50006/2020 DOI 10.54499/UIDB/50006/2020], through national funds; the Consellería de Cultura, Educación e Universidade (Xunta de Galicia) under the scope of the strategic funding of Competitive Reference Group [grant number ED431C 2022/03-GRC], the “Centro singular de investigación de Galicia” (accreditation 2019–2022) funded by the European Regional Development Fund (ERDF) [grant number ED431G2019/06]. The authors also acknowledge the

postdoctoral fellowship of Martín Pérez-Pérez, funded by Xunta de Galicia [grant number ED481D-2023-003]. Funding for open access charge: Universidade de Vigo/CISUG.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.jcs.2024.103960>.

References

- Alfred, J., 2001. Mining the bibliome. *Nat. Rev. Genet.* <https://doi.org/10.1038/35076512>.
- Amirbakhhtar, N., Ismaili, A., Ghaffari, M.R., Mansuri, R.M., Sanjari, S., Shobbar, Z.S., 2021. Transcriptome analysis of bread wheat leaves in response to salt stress. *PLoS One.* <https://doi.org/10.1371/journal.pone.0254189>.
- Amjad Khan, W., Chun-Mei, H., Khan, N., Iqbal, A., Lyu, S.W., Shah, F., 2017. Bioengineered plants can be a useful source of omega-3 fatty acids. *BioMed Res. Int.* <https://doi.org/10.1155/2017/7348919>.
- Araus, J.L., Serret, M.D., Lopes, M.S., 2019. Transgenic solutions to increase yield and stability in wheat: shining hope or flash in the pan? *J. Exp. Bot.* <https://doi.org/10.1093/jxb/erz077>.
- Balcūnaitė-Murzienė, G., Dzikaras, M., 2021. Wheat germ agglutinin—from toxicity to biomedical applications. *Appl. Sci.* <https://doi.org/10.3390/app11020884>.
- Bolland, J.M., 1988. Sorting out centrality: an analysis of the performance of four centrality models in real and simulated networks. *Soc. Network.* 10, 233–253. [https://doi.org/10.1016/0378-8733\(88\)90014-7](https://doi.org/10.1016/0378-8733(88)90014-7).
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., others, 2023. Vicuna: an open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See. <https://vicuna.lmsys.org>. (Accessed 14 April 2023).
- Cuatrecasas, P., 1973. Interaction of concanavalin A and wheat germ agglutinin with the insulin receptor of fat cells and liver. *J. Biol. Chem.* [https://doi.org/10.1016/s0021-9258\(19\)43962-8](https://doi.org/10.1016/s0021-9258(19)43962-8).
- de Sousa, T., Ribeiro, M., Sabeça, C., Igrejas, G., 2021. The 10,000-year success story of wheat. *Foods.* <https://doi.org/10.3390/foods10092124>.
- Frakolaki, G., Giannou, V., Topakas, E., Tzia, C., 2018. Chemical characterization and breadmaking potential of spelt versus wheat flour. *J. Cereal. Sci.* <https://doi.org/10.1016/j.jcs.2017.08.023>.
- Geonames.org, 2016. GeoNames Database [WWW Document]. (Accessed 26 July 2016).
- Geopy, 2023 [WWW Document]. <https://geopy.readthedocs.io/en/stable/index.html>. (Accessed 8 October 2023).
- Giraldo, P., Benavente, E., Manzano-Agugliaro, F., Gimenez, E., 2019. Worldwide research trends on wheat and barley: a bibliometric comparative analysis. *Agronomy.* <https://doi.org/10.3390/agronomy9070352>.
- Gu, Y., Timn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., Poon, H., 2022. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthc.* 3 <https://doi.org/10.1145/3458754>.
- Gupta, P.K., 2024. Drought-tolerant transgenic wheat HB4®: a hope for the future. *Trends Biotechnol.* <https://doi.org/10.1016/j.tibtech.2023.12.007>, 0.
- Jabran, M., Ali, M.A., Zahoor, A., Muhae-Ud-Din, G., Liu, T., Chen, W., Gao, L., 2023. Intelligent reprogramming of wheat for enhancement of fungal and nematode disease resistance using advanced molecular techniques. *Front. Plant Sci.* <https://doi.org/10.3389/fpls.2023.1132699>.
- Jost, M., Outram, M.A., Dibley, K., Zhang, J., Luo, M., Ayliffe, M., 2023. Plant and pathogen genomics: essential approaches for stem rust resistance gene stacks in wheat. *Front. Plant Sci.* <https://doi.org/10.3389/fpls.2023.1223504>.
- Kim, H.S., Patel, K.G., Orosz, E., Kothari, N., Demeyn, M.F., Pyrsopoulos, N., Ahlawat, S. K., 2016. Time trends in the prevalence of celiac disease and gluten-free diet in the US population: results from the national health and nutrition examination surveys 2009–2014. *JAMA Intern. Med.* <https://doi.org/10.1001/jamainternmed.2016.5254>.
- Kupper, C., 2005. Dietary guidelines and implementation for celiac disease. *Gastroenterology.* <https://doi.org/10.1053/j.gastro.2005.02.024>.
- Lebwohl, B., Cao, Y., Zong, G., Hu, F.B., Green, P.H.R., Neugut, A.I., Rimm, E.B., Sampson, L., Dougherty, L.W., Giovannucci, E., Willett, W.C., Sun, Q., Chan, A.T., 2017. Long term gluten consumption in adults without celiac disease and risk of coronary heart disease: prospective cohort study. *BMJ.* <https://doi.org/10.1136/bmj.j1892>.
- Lebwohl, B., Sanders, D.S., Green, P.H.R., 2018. Coeliac disease. *Lancet.* [https://doi.org/10.1016/s0140-6736\(17\)31796-8](https://doi.org/10.1016/s0140-6736(17)31796-8).
- Leetaru, K., 2024. Visualizing an entire day of global news coverage: technical experiments: PCA vs UMAP for HDBSCAN & t-SNE dimensionality reduction – the GDELT project. <https://blog.gdeltproject.org/visualizing-an-entire-day-of-global-news-coverage-technical-experiments-pca-vs-umap-for-hdbscan-t-sne-dimensionality-reduction/>, 2.6.24.
- Li, L., Song, D., Ma, R., Qiu, X., Huang, X., 2021. KNN-BERT: fine-tuning pre-trained models with KNN classifier. <https://doi.org/10.48550/arxiv.2110.02523>.
- Lyson, H.C., Le, G.M., Zhang, J., Rivadeneira, N., Lyles, C., Radcliffe, K., Pasick, R.J., Sawaya, G., Sarkar, U., Centola, D., 2019. Social media as a tool to promote health awareness: results from an online cervical cancer prevention study. *J. Cancer Educ.* 34, 819–822. <https://doi.org/10.1007/s13187-018-1379-8>.

- Ma, H., 2023. Editorial for the special issue “genetics studies on wheat.”. *Genes*. <https://doi.org/10.3390/genes14091761>.
- Makwana, P., Kodinariya, T.M., Makwana, P.R., 2013. Review on determining of cluster in K-means clustering review on determining number of cluster in K-means clustering. *Int. J. Adv. Res. Comput. Sci. Manag. Stud.* 1.
- McInnes, L., Healy, J., Saul, N., Großberger, L., 2018. UMAP: Uniform Manifold approximation and projection. *J. Open Source Softw.* <https://doi.org/10.21105/joss.00861>.
- Nieminen, J., 1974. On the centrality in a graph. *Scand. J. Psychol.* 15, 332–336. <https://doi.org/10.1111/j.1467-9450.1974.tb00598.x>.
- Qin, X., Zhang, F., Liu, C., Yu, H., Cao, B., Tian, S., Liao, Y., Siddique, K.H.M., 2015. Wheat yield improvements in China: past trends and future directions. *Field Crops Res.* <https://doi.org/10.1016/j.fcr.2015.03.013>.
- Ribeiro, M., Nunes-Miranda, J.D., Branlard, G., Carrillo, J.M., Rodríguez-Quijano, M., Igrejas, G., 2013. One hundred years of grain omics: identifying the glutens that feed the world. *J. Proteome Res.* <https://doi.org/10.1021/pr400663t>.
- Ribeiro, M., Nunes, F.M., 2019. We might have got it wrong: modern wheat is not more toxic for celiac patients. *Food Chem.* <https://doi.org/10.1016/j.foodchem.2018.12.003>.
- Ribeiro, M., Picascia, S., Rhazi, L., Gianfrani, C., Carrillo, J.M., Rodríguez-Quijano, M., Branlard, G., Nunes, F.M., 2019. Effect of in situ gluten-chitosan interlocked self-assembled supramolecular architecture on rheological properties and functionality of reduced celiac-toxicity wheat flour. *Food Hydrocolloids.* <https://doi.org/10.1016/j.foodhyd.2018.12.026>.
- Sabença, C., Ribeiro, M., de Sousa, T., Poeta, P., Bagulho, A.S., Igrejas, G., 2021. Wheat/gluten-related disorders and gluten-free diet misconceptions: a review. *Foods.* <https://doi.org/10.3390/foods10081765>.
- Salisbury, L., 2020. Scopus CiteScore and Clarivate journal citation reports. *Charlest. Advis.* <https://doi.org/10.5260/chara.21.4.5>.
- Searls, D.B., 2001. Mining the bibliome. *Pharmacogenomics J.* <https://doi.org/10.1038/sj.tpj.6500030>.
- Shewry, P.R., 2009. Wheat. *J. Exp. Bot.* <https://doi.org/10.1093/jxb/erp058>.
- Shewry, P.R., Ward, J.L., 2012. Exploiting genetic variation to improve wheat composition for the prevention of chronic diseases. *Food Energy Secur.* (2) <https://doi.org/10.1002/fes3>.
- United Nations- DESA, 2022. *World Population Prospects 2022*, United Nation.
- Wieser, H., Koehler, P., Scherf, K.A., 2023. Chemistry of wheat gluten proteins: quantitative composition. *Cereal Chem.* <https://doi.org/10.1002/cche.10553>.
- Wood, T.A., 2022. Country named entity recognition (computer software) [WWW Document]. URL, Version 0.4. <https://fastdatascience.com/country-named-entity-recognition/>. (Accessed 8 October 2023).
- Worldometer, 2022. *World Population Clock: 8 Billion People (LIVE, 2022)* - Worldometer [WWW Document]. Worldometer.
- Yang, W., Shen, Y., 2018. Quality assessment of feed wheat in ruminant diets. In: *Global Wheat Production.* <https://doi.org/10.5772/intechopen.75588>.