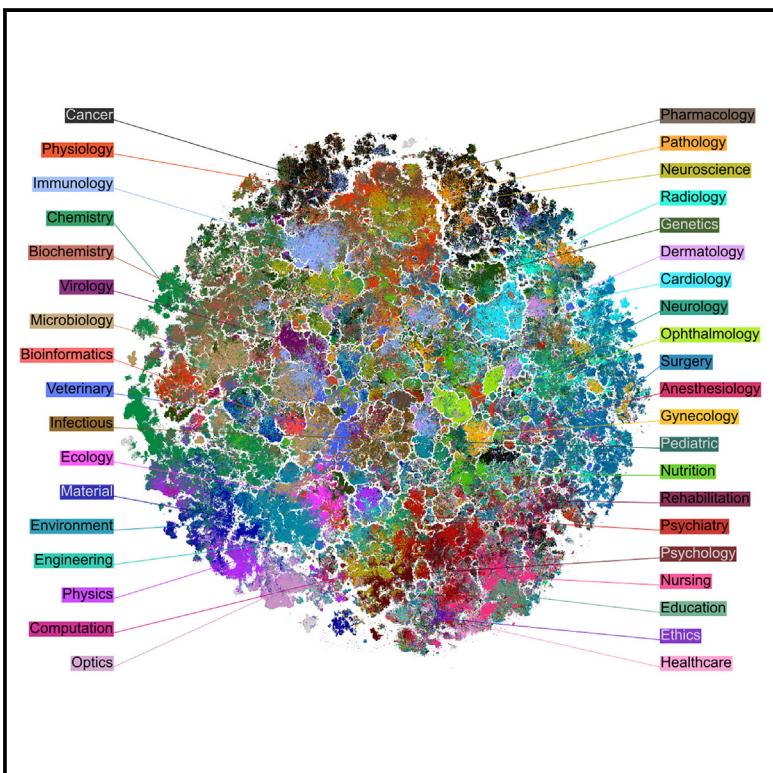


Graphical abstract



Highlights

- We develop a 2D map of biomedical papers based on abstract texts
- Our embedding contains 21 million papers covering all areas of biology and medicine
- The map can be used to study issues ranging from gender bias to fraudulent research

Article

Authors

Rita González-Márquez, Luca Schmidt, Benjamin M. Schmidt, Philipp Berens, Dmitry Kobak

Correspondence

dmitry.kobak@uni-tuebingen.de

In brief

This study presents a 2D map based on the abstracts of biomedical research articles from PubMed. Containing 21 million English articles, this map highlights several publishing issues, including gender bias and fraudulent research.



Article

The landscape of biomedical research

Rita González-Márquez,^{1,2} Luca Schmidt,^{1,2} Benjamin M. Schmidt,³ Philipp Berens,^{1,2} and Dmitry Kobak^{1,2,4,5,*}

¹Hertie Institute for AI in Brain Health, University of Tübingen, Germany

²Tübingen AI Center, Tübingen, Germany

³Nomic AI, New York, New York, USA

⁴IWR, Heidelberg University, Heidelberg, Germany

⁵Lead contact

*Correspondence: dmitry.kobak@uni-tuebingen.de

<https://doi.org/10.1016/j.patter.2024.100968>

THE BIGGER PICTURE Over 1.5 million scientific articles on biomedicine and life sciences are now published and collected in the PubMed database every year. This vast scale makes it challenging to see how biomedicine evolves in time. Large language models can produce embeddings of huge text corpora and can thereby be leveraged to provide innovative visualizations of the scientific literature, as shown in this work.

SUMMARY

The number of publications in biomedicine and life sciences has grown so much that it is difficult to keep track of new scientific works and to have an overview of the evolution of the field as a whole. Here, we present a two-dimensional (2D) map of the entire corpus of biomedical literature, based on the abstract texts of 21 million English articles from the PubMed database. To embed the abstracts into 2D, we used the large language model PubMedBERT, combined with t-SNE tailored to handle samples of this size. We used our map to study the emergence of the COVID-19 literature, the evolution of the neuroscience discipline, the uptake of machine learning, the distribution of gender imbalance in academic authorship, and the distribution of retracted paper mill articles. Furthermore, we present an interactive website that allows easy exploration and will enable further insights and facilitate future research.

INTRODUCTION

The rate of scientific publishing has been increasing constantly over the past century,^{1,2} with over 1 million articles being currently published every year in biomedicine and life sciences alone. Information about academic publications in these fields is collected in the PubMed database, maintained by the United States National Library of Medicine (pubmed.ncbi.nlm.nih.gov). It now contains over 35 million scientific papers from the last 50 years.

This rapid growth of the biomedical literature makes it difficult to track the evolution of biomedical publishing as a whole. Search engines such as PubMed and Google Scholar allow researchers to find specific papers given suitable keywords and to follow the citation networks that these papers are embedded in, yet none of them allows exploration of the biomedical literature landscape from a global perspective. This makes it hard to see how research topics evolve over time, how different fields are related to each other, or how new methods and techniques are adopted in different fields. What is needed to answer such questions is a bird's eye view on the biomedical literature.

In this work we develop an approach that enables all of the above: a global two-dimensional (2D) atlas of the biomedical and life science literature that is based on the abstracts of all 21 million English language articles contained in the PubMed database. For simplicity, our map is based on the abstract texts alone, and does not rely on other article parts, such as its main text, figures, or references. To create the map, we embedded the abstracts into two dimensions using the transformer-based large language model PubMedBERT³ combined with the neighbor-embedding method t-SNE,⁴ adapted to handle samples of this size. Our approach allowed us to create a map with the level of detail substantially exceeding previous works.^{5,6}

We argue that our visualization facilitates exploration of the biomedical literature and can reveal aspects of the data that would not be easily noticed with other analysis methods. We showcase the power of our approach in five examples: we studied (1) the emergence of the COVID-19 literature, (2) the evolution of different subfields of neuroscience, (3) the uptake of machine learning in the life sciences, (4) the distribution of gender imbalance across biomedical fields, and (5) the distribution of retracted paper mill articles. In all cases, we used the embedding



Table 1. Quality metrics for the embeddings

Data	Dim.	Acc. (%)	RMSE	Recall (%)
PubMedBERT	768	69.7	8.4	–
TF-IDF	4,679,130	65.2	8.8	–
t-SNE(BERT)	2	62.6	10.2	6.2
t-SNE(TF-IDF)	2	50.6	11.2	0.7
Chance	–	4.3	12.4	0.0

Acc., kNN accuracy ($k = 10$) of label prediction; RMSE, root mean-squared error of kNN prediction of publication year; Recall, overlap between k nearest neighbors in the 2D embedding and in the high-dimensional space. See [experimental procedures](#) for details.

to formulate specific hypotheses about the data that were later confirmed by a dedicated statistical analysis of the original high-dimensional dataset.

The resulting map of the biomedical research landscape is publicly available as an interactive web page at <https://static.nomic.ai/pubmed.html>, developed using the deepscatter library.⁷ It allows users to navigate the atlas, zoom, and search by article title, journal, and author names, while loading individual scatter points on demand. We envisage that the interactive map will allow further insights into the biomedical literature, beyond the ones we present in this work.

RESULTS

2D atlas allows to explore the PubMed database

We downloaded the complete PubMed database (2021 snapshot) and, after initial filtering (see [experimental procedures](#)), were left with 20,687,150 papers with valid English abstracts, the majority of which (99.8%) were published in 1970–2021 (Figure S1). Our goal was to generate a 2D embedding of the abstract texts to facilitate exploration of the data.

To annotate our atlas, we chose a set of 38 labels covering basic life science fields such as “virology” and “biochemistry,” and medical specialties such as “radiology” and “ophthalmology.” We assigned each label to the papers published in journals with the corresponding word in journal titles. For example, all papers published in *Annals of Surgery* were labeled “surgery.” As a result, 34.4% of all papers received a label, while the rest remained unlabeled. This method misses papers published in interdisciplinary journals such as *Science* or *Nature*, but labels the core works in each discipline. We chose our labels so that they would cover every region of the 2D space. Therefore, despite only having 34% of the papers labeled, we consider this fraction to be representative of the whole landscape.

To generate a 2D map of the entire PubMed database, we first obtained a 768-dimensional numerical representation of each abstract using PubMedBERT,³ which is a Transformer-based⁸ language model trained on PubMed abstracts and full-text articles from PubMed Central. We then reduced the dimensionality to two using t-SNE.⁴

For the initial step of computing a numerical representation of the abstracts, we evaluated several text processing methods, including bag-of-words representations such as TF-IDF (term frequency-inverse document frequency)⁹ and several other BERT-derived models, including the original BERT,¹⁰ SBERT,¹¹

SciBERT,¹² BioBERT,¹³ SPECTER,¹⁴ SimCSE,¹⁵ and SciNCL.¹⁶ We chose PubMedBERT because it best grouped papers together in terms of their label, quantified by the k nearest neighbor (kNN) classification accuracy when each label is predicted based on the most frequent label of its 10 nearest neighbors (Table S2). For the PubMedBERT representation, this prediction was correct 69.7% of the time (Table 1). For comparison, TF-IDF, which is simpler and faster to compute, yielded lower kNN accuracy (65.2%).

For the second step, we used t-SNE with several modifications that allowed us to run it effectively on very large datasets. These modifications included uniform affinities to reduce memory consumption and extended optimization to ensure better convergence (see [experimental procedures](#)). With these modifications, t-SNE performs better than other neighbor-embedding methods such as UMAP¹⁷ in terms of kNN accuracy and memory requirements.¹⁸ The resulting embedding showed good label separation, with kNN accuracy in 2D of 62.6%, not much worse than the 4,679,130-dimensional TF-IDF representation.

We interpret the resulting embedding as the map of the biomedical literature (Figure 1). It showed sensible global organization, with natural sciences mainly located on the left side and medical specialties gathered on the right side; physics- and engineering-related works occupied the bottom-left part (Figures S2 and S3). Related disciplines were located next to each other: for example, the biochemistry region was overlapping with chemistry, whereas psychology was merging into psychiatry. A t-SNE embedding based on the TF-IDF representation had similar large-scale structure but worse kNN accuracy (50.6%; Figure S4).

In addition to this global structure, the map revealed rich and detailed fine structure and was fragmented into small clusters containing hundreds to thousands of papers each (Figure S5A). Even though immediate neighborhoods were distorted compared with the 768-dimensional PubMedBERT representation (only 6.2% of the nearest neighbors in \mathbb{R}^2 were nearest neighbors in \mathbb{R}^{768} ; we call this metric kNN recall), manual inspection of the clusters suggested that they consisted of papers on clearly defined narrow topics.

Moreover, the map had rich temporal structure, with papers of the same age tending to be grouped together (Figure 1B). While this structure may be influenced by changes in writing style and common vocabulary, it is likely primarily caused by research topics evolving over time and becoming more or less fashionable. The most striking example of this effect is a cluster of very recent papers published in 2020–2021 that is very visible in the middle of the map (bright yellow in Figure 1B). We will use this island as our first example of how the map can be used to guide understanding of the publishing landscape and how it allows to form hypotheses about the structure and temporal evolution of biomedical research. We will show that these hypotheses can be rigorously confirmed in the high-dimensional embedding space.

The COVID-19 literature is uniquely isolated

The bright yellow island we identified above comprised works related to COVID-19 (Figure 1B), with 85% of papers on COVID-related topics, and 15% on other respiratory epidemics. Our dataset included in total 132,802 COVID-related papers

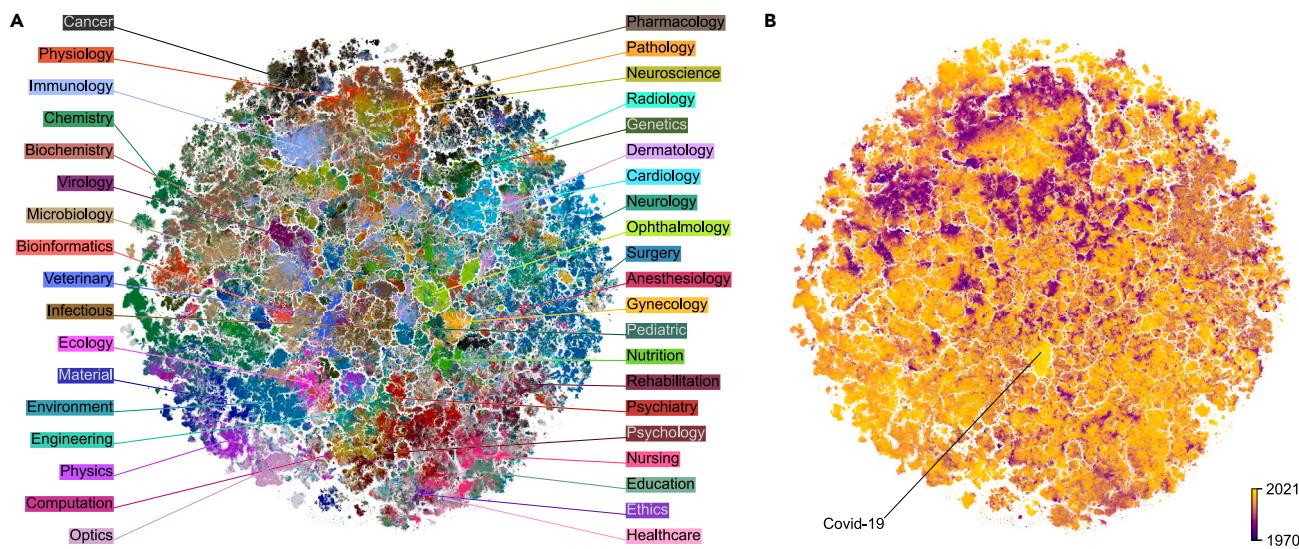


Figure 1. 2D embedding of the PubMed dataset

Paper abstracts ($n = 21$ million) were transformed into 768-dimensional vectors with PubMedBERT³ and then embedded in 2D with t-SNE.⁴ (A) Colored using labels based on journal titles. Unlabeled papers are shown in gray and are displayed in the background. (B) Colored by publication year (dark, 1970 and earlier; light, 2021).

(based on terms such as COVID-19, SARS-CoV-2, etc., present in their abstracts; see [experimental procedures](#)), which constituted 5.2% of all PubMed papers published in 2020–2021. As the pandemic and its effects were studied by many different biomedical fields, one might have expected the COVID papers to be distributed across the embedding in their corresponding disciplines. Instead, most (59.3%) of the COVID-related papers were grouped together in one cluster, while the rest were sparsely distributed across the map (Figure S6A).

The main COVID cluster was surrounded by articles on other epidemics, public health issues, and respiratory diseases. When we zoomed in, we found rich inner structure within the COVID cluster itself, with multiple COVID-related topics separated from each other (Figure 2). Papers on mental health and societal impact, on public health and epidemiological control, on immunology and vaccines, on clinical symptoms and treatment were all largely non-overlapping, and were further divided into even narrower sub-fields. This suggests that our map can be useful for navigating the literature on the scale of narrow and focused scientific topics.

Seeing that the COVID papers prominently stood out in the map (Figure 1B), we hypothesized that the COVID literature was more isolated from the rest of the biomedical literature, compared with other similar fields. To test this, we selected several comparable sets of papers, such as papers on HIV/AIDS or influenza, or all papers published in virology or ophthalmology journals (two labels that appeared particularly compact in Figure 1A). We measured the *isolatedness* of each corpus in the high-dimensional space by the fraction of their k NNs that belonged to the same corpus. We found that, indeed, COVID literature had the highest isolatedness, in both BERT (80.6%) and TF-IDF (76.2%) representations (Table 2). This suggests that the COVID-19 pandemic had an unprecedented effect on the scientific literature, creating a separate and uniquely detached field of study in only 2 years.

We investigated the driving factors behind the emergence of the COVID island in the 2D space using the TF-IDF representation and saw that, even though the presence of COVID keywords (such as “COVID” or “SARS-CoV”) did play some role in the island formation, it was not the only source of similarity between COVID papers (Figure S7).

Changing focus within neuroscience

As we have seen in the extreme example of the COVID literature, the atlas can be used to study composition and temporal trends across disciplines. We next show how it can also provide insights into shifting topics and trends inside a discipline. We demonstrate this using the example of neuroscience. Neuroscience papers ($n = 240,135$) in the map were divided into two main clusters (Figure 3A). The upper one contained papers on molecular and cellular neuroscience, while the lower one consisted of studies on behavioral and cognitive neuroscience. While it has been shown that articles in brain-related journals show a separation between basic science and clinical applications,¹⁹ our map revealed a different bimodality, separating cellular from behavioral neuroscience. Several smaller clusters comprised papers on neurodegenerative diseases and sensory systems.

Coloring this part of the embedding by publication year indicated that the cellular/molecular region on average had older papers than the cognitive/behavioral region (Figure 3B). This suggests that the relative publication volume in different subfields of neuroscience has changed with time. To test this hypothesis directly, we devised a metric measuring the overlap between neuroscience and any given related discipline across time. We defined k NN overlap as the fraction of k NNs of neuroscience papers that belonged to a given discipline in the high-dimensional space. We found that the overlap of neuroscience with physiology and pharmacology has decreased since the 1970s, while its overlap with psychiatry, psychology, and computation has

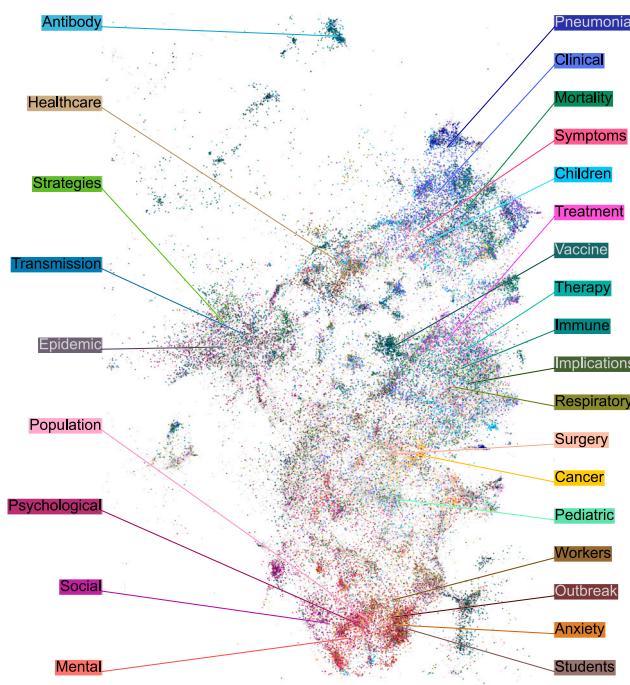


Figure 2. COVID-19 region of the map

Colors are assigned using the most common keywords appearing in paper titles. Uncolored COVID papers are shown in the background in gray. This region in the embedding also contained some non-COVID papers (~ 15%) about other respiratory epidemics; they are not shown.

increased, in particular after the 1990s (Figure 3C). Indeed, neuroscience originated as a study of the nervous system within physiology, but gradually broadened its scope to include cognitive neuroscience, related to psychology, as well as computational neuroscience, related to computer science and machine learning.

The uptake of machine learning

In recent years, computational methods and machine learning have increasingly found use in various biomedical disciplines.²⁰ To explore the use of machine learning (ML) in the biomedical landscape, we computed the fraction of papers claiming to use machine learning (defined as papers mentioning “machine learning” in their abstracts) within different medical disciplines across time (Figure 4A). We found that the uptake of ML differed substantially across disciplines. Radiology was the first discipline to show an increase in ML adoption, shortly after 2015, followed by psychiatry and neurology. In oncology, ML adoption started later but showed accelerated rise over the last 5 years. This is in contrast with specialties such as dermatology and gynecology, which did not see any ML usage until ~2020.

However, this simple analysis is constrained by the set of labels chosen *a priori*. The 2D map allows an unbiased exploration that does not rely on labels. For that, we highlighted all machine learning papers ($n = 38,446$) in the embedding (Figure 4B). Papers claiming to use machine learning were grouped in the map into several clusters, covering topics ranging from computational biology to healthcare data management. These ML papers were more prevalent in the life science half of the map (left) and rather

rare in the medical part (right). Within the medical part of the corpus, ML papers were concentrated in several regions, such as analysis of tumor imaging (radiology) or cancer biomarkers (oncology).

To further explore the ML-heavy regions, we selected and manually labeled 12 of them (Figure 4B) and computed the fraction of papers mentioning specific ML and statistical methods (Table S1). We found that the usage of ML techniques varied strongly across regions. Deep learning and convolutional networks were prominent in the image segmentation region (with applications, e.g., in microscopy). Clustering was often used in analyzing sequencing data. Neural networks and support vector machines were actively used in structural biology. Principal component analysis was important for data analysis in mass spectrometry.

We expanded this analysis to the whole corpus by identifying 342,070 papers (1.7%) mentioning the same ML and statistical methods in their abstracts (Figure 4C). We found that the medical part of the embedding was dominated by classical linear methods such as linear regression and factor analysis, whereas more modern nonlinear and nonparametric methods were mostly used in non-medical research. This shows that the medical disciplines are being slower in taking up new computational techniques compared with basic life sciences.

Exploring the gender gap

In this section we show how the map can be used to explore and better understand social disparities in biomedical publishing such as the extent and distribution of the well-known gender imbalance in academic authorship.^{21–25} We used the first name (where available) of the first and the last author of every PubMed paper to infer their gender using the gender tool.²⁶ The gender inference is only approximate, as many first names were absent in the US-based training data, biasing our analysis toward Western academia, and some names are inherently gender-ambiguous (see *experimental procedures*). Overall, this procedure allowed us to infer the gender of 62.3%/63.1% first/last authors with available first names. Among those, 42.4% of first authors and 29.1% of last authors were female. While some academic fields, such as mathematics and physics, tend to prefer alphabetic ordering of the authors, in biomedicine the first author is usually the trainee (PhD student or postdoc) who did the practical hands-on project work and the last author is the supervisor or principal investigator.

The fraction of female authors steadily increased with time (Figure 5A), with first and last authors being 47.2% and 34.4% female in 2021. We found a delay of ~20 years between the first and the last author curves, suggesting that it takes more than one academic generation for the differences in gender bias to propagate from mentees to mentors.

Within most individual disciplines, the fraction of female first authors increased with time (Figure 5B), even in disciplines where this fraction was already high, such as education (increased from 55% female in 2005 to 60% in 2020). This increase also happened in male-dominated fields such as computation, physics, or surgery (increase from 15% to 25%). Notably, the female proportion in material sciences showed only a modest increase, while nursing, the most female-dominated discipline across all our labels (80.4%), showed a moderate decrease.

Table 2. Isolatedness metric for several sets of papers

	n	BERT (%)	TF-IDF (%)
COVID-19	132,802	80.6	76.2
HIV/AIDS	308,077	63.9	62.3
Influenza	90,575	57.9	64.1
Meta-analysis	145,358	52.6	38.5
Virology	112,807	47.7	39.1
Ophthalmology	144,411	47.7	43.6

Fraction of k nearest neighbors of papers from each corpus that also belong to the same corpus (see [experimental procedures](#)). The first four rows show corpora selected based on the abstract text; the last two, based on the journal name.

Our map, when colored by gender, also showed that female authors were not equally distributed across the biomedical publishing landscape ([Figure 5C](#)). First and last female authors were most frequent in the lower right corner of the embedding, covering fields such as nursing, education, and psychology. Furthermore, the map allowed us to explore gender bias beyond the discipline level, revealing a substantial heterogeneity of gender ratios within individual disciplines. For example, in healthcare (overall 49.6% female first authors), there were male- and female-dominated regions in the map. One of the more male-dominated clusters (33.9% female) focused on financial management, while one of the more female ones (68.1% female) focused on patient care ([Figure 5C](#)). In education (58.6% female authors), female authors dominated research on nursing training, whereas male authors were more frequent in research on medical training ([Figure 5D](#)). In surgery, only 24.4% of the first authors were female, but this fraction increased to 61.1% in the cluster of papers on veterinary surgery ([Figure 5E](#)). This agrees with veterinary medicine being a predominantly female discipline (52.2% in total, [Figure 5G](#)). Importantly, these details are lost when averaging across a *priori* labels, while the embedding can suggest the relevant level of granularity.

Retracted papers highlight suspicious literature

We identified 11,756 papers flagged as retracted by PubMed and still having intact abstracts (not containing words such as “retracted” or “withdrawn”; see [experimental procedures](#)). These papers were not distributed uniformly over the 2D map ([Figure 6](#)) but instead concentrated in several specific areas, in particular on top of the map, covering research on cancer-related drugs, marker genes, and microRNA. These areas are known targets of paper mills,^{27–29} which are organizations that produce fraudulent research papers for sale.

Our map is based solely on textual similarity between abstracts. This suggests that non-retracted papers from the regions with high concentration of retracted papers may require an investigation, as their abstracts are similar to the ones from paper mill products. As an example, we considered a region with particularly high fraction (45/422) of retracted papers (second inset in [Figure 6](#)) and randomly selected 25 non-retracted papers for manual inspection. They had similar title format (variations of “MicroRNA-X does Y by targeting Z in osteosarcoma”³⁰), paper structure, and figure style, and 24/25 of them

had authors affiliated with Chinese hospitals—features that are often shared by paper mill products.^{29,31–36} Moreover, many areas with high fraction of retractions consisted of papers stemming mostly from a single country, typically China ([Figure S8](#)), which could by itself be an indicator of paper mill activity.

After we conducted our analysis, the Retraction Watch database of retracted papers was made open to the public. Using their database, we identified an additional 3,572 papers in our map that were not marked as retracted in PubMed but were in fact retracted (red dots in [Figure 6](#)). They were mostly located in the same areas of the map that we identified as suspicious above, validating our conclusions. This does not guarantee that all papers in these areas are fraudulent, but confirms that our 2D map can be used to highlight papers requiring further editorial investigation.³⁷ If additional paper mills are discovered in the future, our map will help to highlight literature clusters requiring further scrutiny.

DISCUSSION

We developed a 2D atlas of the biomedical literature based on the PubMed collection of 21 million paper abstracts using a transformer-based language model (PubMedBERT) and a neighbor-embedding visualization (*t*-SNE) tailored to handle large document libraries. We used this atlas as an exploration tool to study the biomedical research landscape, generating hypotheses that we later confirmed using the original high-dimensional data. Using five distinct examples—the emergence of the COVID-19 literature, the evolution of the neuroscience discipline, the uptake of machine learning, the gender imbalance, and the concentration of retracted fraudulent papers—we argued that 2D visualizations of text corpora can help uncover aspects of the data that other analysis methods may fail to reveal.

We also developed an interactive web version of the embedding (<https://static.nomic.ai/pubmed.html>) based on the deepscatter library,⁷ which allows to navigate the atlas, zoom, and search by title, journal, or author names. In deepscatter, individual points are loaded on demand when zooming-in, like when navigating geographical maps in the browser. This interactive website contains a separate embedding of the latest PubMed data, including 2022–2023 papers ([Figures S9–S10](#)). We plan on updating the visualization in the future using annual PubMed releases.

Neighbor-embedding methods such as *t*-SNE have known limitations. For the datasets of our size, the few closest neighbors in the 2D embedding space are typically different from the neighbors in the high-dimensional BERT representation ([Table 1](#)). This makes our map suboptimal for finding the most similar papers to a given query paper, and other tools, such as conventional (Google Scholar, PubMed) or citation-based ([connectedpapers.com](#)) search engines, may be more appropriate for this task. Instead, our map is useful for navigating the literature on the scale of narrow and focused scientific topics. Neighbor-embedding algorithms can misrepresent the global organization of the data.^{38–40} We used methods designed to mitigate this issue^{18,39,41} and, indeed, found that related research areas were located close to each other.

In annotating our atlas, we selected 38 labels spanning various life science fields and medical specialties. Each label

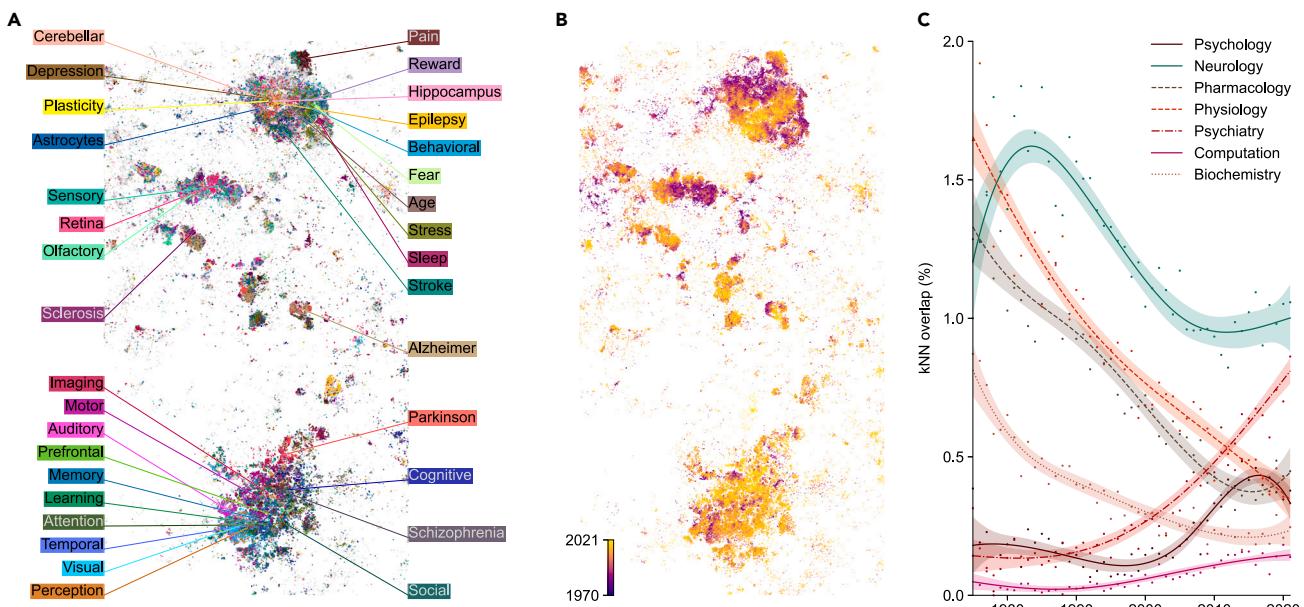


Figure 3. Neuroscience literature

(A) Articles published in neuroscience journals, colored by presence of specific keywords in paper titles.

(B) The same articles colored by the publication year (dark, 1970 and earlier; light, 2021).

(C) Fraction of the high-dimensional kNNs of neuroscience papers that belonged to a given discipline (biochemistry, computation, neurology, pharmacology, physiology, psychiatry, psychology). We chose to analyze those disciplines because they had the highest confusion scores with the neuroscience class in a kNN classifier. Points: yearly averages. Smooth curves and 95% confidence intervals were obtained with generalized additive models (see [experimental procedures](#)).

was assigned to papers in journals with the corresponding word in their titles, resulting in 34.4% of papers being labeled. Although this method overlooks interdisciplinary journals such as *Science* or *Nature*, it ensures that core works in each discipline are labeled. While PubMed uses several systems to organize articles into categories (such as keywords or Medical Subject Headlines [MeSH] terms), creating labels based on them would likely require a more involved manual curation process. We found that our labels covered most of the 2D space ([Figure 1](#)), and therefore considered the labeled subset representative of the entire landscape.

Our atlas provides the most detailed visualization of the biomedical literature landscape to date. Previously, PubMed abstracts were clustered based on textual bag-of-words similarity and citation information, and the clusters were displayed using a 2D embedding.⁵ Their map exhibits similar large-scale organization, but only shows 29,000 clusters, so our map is almost three orders of magnitude more detailed. The BioBERT model was previously applied to the PubMed dataset to extract information on biomedical concepts, such as proteins or drugs.⁴² Previous work on visualizing large text corpora includes Schmidt⁴³ and González-Márquez et al.¹⁸ Both were based on bag-of-words representations of the data. Here, we showed that BERT-based models outperform TF-IDF for representing scientific abstracts.

An alternative approach to visualizing collections of academic works is to use information on citations as a measure of similarity, as opposed to semantic or textual similarity. For example, [paperscape.org](#) visualizes 2.2 million papers from the arXiv preprint server using a force-directed layout of the citation graph.

Similarly, [opensyllabus.org](#) uses node2vec⁴⁴ and UMAP to visualize 1.1 million texts based on their co-appearance in the US college syllabi. Similar approach was used by Noichl⁴⁵ to visualize 68,000 articles on philosophy based on their reference lists. Here, we based our embedding on the abstract texts alone, and in future work it would be interesting to combine textual and co-citation similarity in one map (citation graph for PubMed papers can be obtained from OpenAlex⁴⁶, MAG⁴⁷ and/or PubMed itself). The functionality of our interactive web version is similar to [opensyllabus.org](#) and [paperscape.org](#), but we successfully display one order of magnitude more points.

We achieved the best representation of the PubMed abstracts using the PubMedBERT model. As the progress in the field of language models is currently very fast, it is likely that a better representation may soon become available. One promising approach could be to train sentence-level models such as SBERT¹¹ on the biomedical text corpus. Another active avenue of research is fine-tuning BERT models using contrastive learning^{15,48} and/or using citation graphs.^{14,16} While we found that these models were outperformed by PubMedBERT, similar methods⁴⁹ could be used to fine-tune the PubMedBERT model itself, potentially improving its representation quality further. Finally, larger generative language models such as recently developed BioGPT⁵⁰ or BioMedLM⁵¹ can possibly lead to better representations as well.

In conclusion, we suggested a novel approach for visualizing large document libraries and demonstrated that it can facilitate data exploration and help generate novel insights. Many further meta-scientific questions can be investigated in the future using our approach.

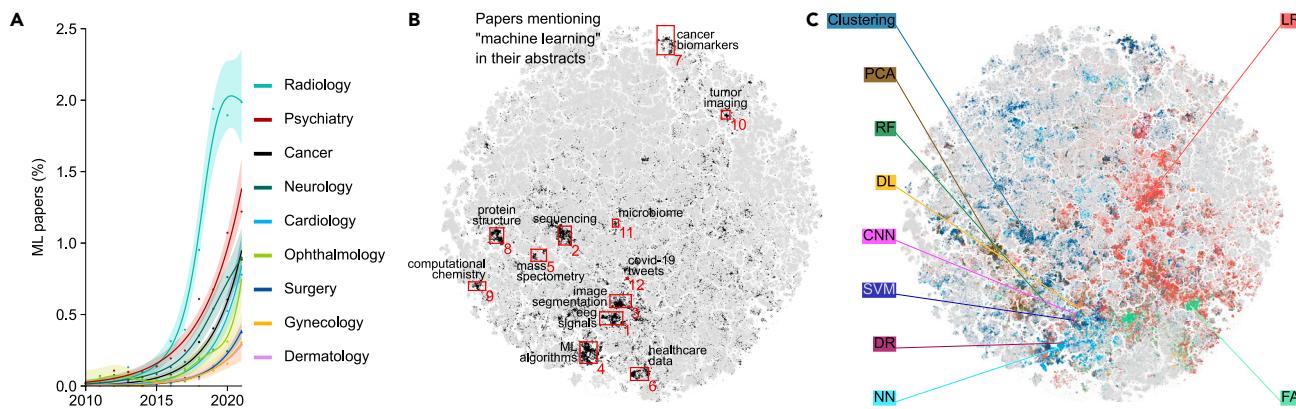


Figure 4. Machine learning papers

(A) Percentage of papers mentioning “machine learning” in their abstracts across time for different medical disciplines. Smooth curves and 95% confidence intervals were obtained using generalized additive models and the points correspond to yearly percentages (see [experimental procedures](#)).
(B) Papers mentioning “machine learning” in their abstracts, grouped into 12 clusters that we manually labeled.
(C) Papers colored according to various statistical and machine learning methods mentioned in their abstracts. PCA, principal component analysis; RF, random forest; DL, deep learning; CNN, convolutional neural network; SVM, support vector machine; DR, dimensionality reduction; NN, neural networks; LR, linear regression; FA, factor analysis. Some of the highlighted NN papers may refer to biological neural networks.

EXPERIMENTAL PROCEDURES

Resource availability

Lead contact

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Dmitry Kobak (dmitry.kobak@uni-tuebingen.de).

Materials availability

This study did not generate new unique reagents.

Data and code availability

- The analysis code is available at <https://github.com/berenslab/pubmed-landscape>. All original code has been deposited at Zenodo under <https://doi.org/10.5281/zenodo.10727578>,⁵² and is publicly available as of the date of publication.
- This paper analyzes existing, publicly available data. It can be obtained by directly accessing the bulk download service (www.nlm.nih.gov/databases/download/pubmed_medicline.html) from PubMed. We made publicly available a processed version of our dataset: a csv.zip file (20,687,150 papers, 1.3 GB) including PMID, title, journal name, publication year, embedding x and y coordinates, our label, and our color used in [Figure 1A](#). We also included two additional files: the raw abstracts (csv.zip file, 9.5 GB) and the 768-dimensional PubMedBERT embeddings of the abstracts (NumPy array in float16 precision, 31.8 GB). They can all be downloaded from Zenodo under <https://doi.org/10.5281/zenodo.7695389>.⁵³
- Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

PubMed dataset

We downloaded the complete PubMed database (295 GB) as XML files using the bulk download service (www.nlm.nih.gov/databases/download/pubmed_medicline.html). PubMed releases a new snapshot of their database every year; they call it a “baseline.” In our previous work¹⁸ we used the 2020 baseline (files called pubmed21n0001.xml.gz to 1062.xml.gz, download date: 26.01.2021). In this work, we supplemented them with the additional files from the 2021 baseline (files called pubmed22n1062.xml.gz to 1114.xml.gz, download date: 27.04.2022). After the analysis was completed, we realized that our dataset had 0.07% duplicate papers; they should not have had any noticeable influence on the reported results.

We used the Python xml package to extract PubMed ID, title, abstract, language, journal title, ISSN, publication date, and author names of all 33.4

million papers. We filtered out all 4.7 million non-English papers, 10.8 million papers with empty abstracts, 0.3 million papers with abstracts shorter than 250 or longer than 4,000 symbols ([Figures S1](#) and [S11](#)), and 27,000 papers with unfinished abstracts. Papers with unfinished abstracts needed to be excluded because otherwise they were grouped together in the BERT representation, creating artifact clusters in the embedding. We defined unfinished abstracts as abstracts not ending with a period, a question mark, or an exclamation mark. Some abstracts ended with a phrase “(ABSTRACT TRUNCATED AT ... WORDS)” with a specific number instead of “...”. We removed all such phrases and analyzed the remaining abstracts as usual, even though they did not contain the entire text of the original abstracts. In some cases, abstracts were divided in subsections (such as methods, results, etc.). We excluded subsection titles so that the resulting abstract had effectively a single paragraph. Overall, we were left with 20,687,150 papers for further analysis.

This collection contains papers from the years 1808–2022. MEDLINE, the largest component of PubMed, started its record in 1966 and later included some noteworthy earlier papers. Therefore, the majority (99.8%) of the PubMed papers are post-1970 ([Figure S1C](#)). There are only few papers from 2022 in our dataset. The 2021 data in this PubMed snapshot were also incomplete.

Label assignment

We labeled the dataset by selecting 38 keywords contained in journal titles that reflected the general topic of the paper. We based our choice of keywords on lists of medical specialties and life science branches that appeared frequently in the journal titles in our dataset. The 38 terms are: anesthesiology, biochemistry, bioinformatics, cancer, cardiology, chemistry, computation, dermatology, ecology, education, engineering, environment, ethics, genetics, gynecology, healthcare, immunology, infectious, material, microbiology, neurology, neuroscience, nursing, nutrition, ophthalmology, optics, pathology, pediatric, pharmacology, physics, physiology, psychiatry, psychology, radiology, rehabilitation, surgery, veterinary, and virology.

Papers were assigned a label if their journal title contained that term, either capitalized or not, and were left unlabeled otherwise. Journal titles containing more than one term were assigned randomly to one of them. This resulted in 7,123,706 labeled papers (34.4%).

Our journal-based labels do not constitute the ground truth for the topic of each paper, and so the highest possible classification accuracy is likely well below 100%. Nevertheless, we reasoned that the higher the classification

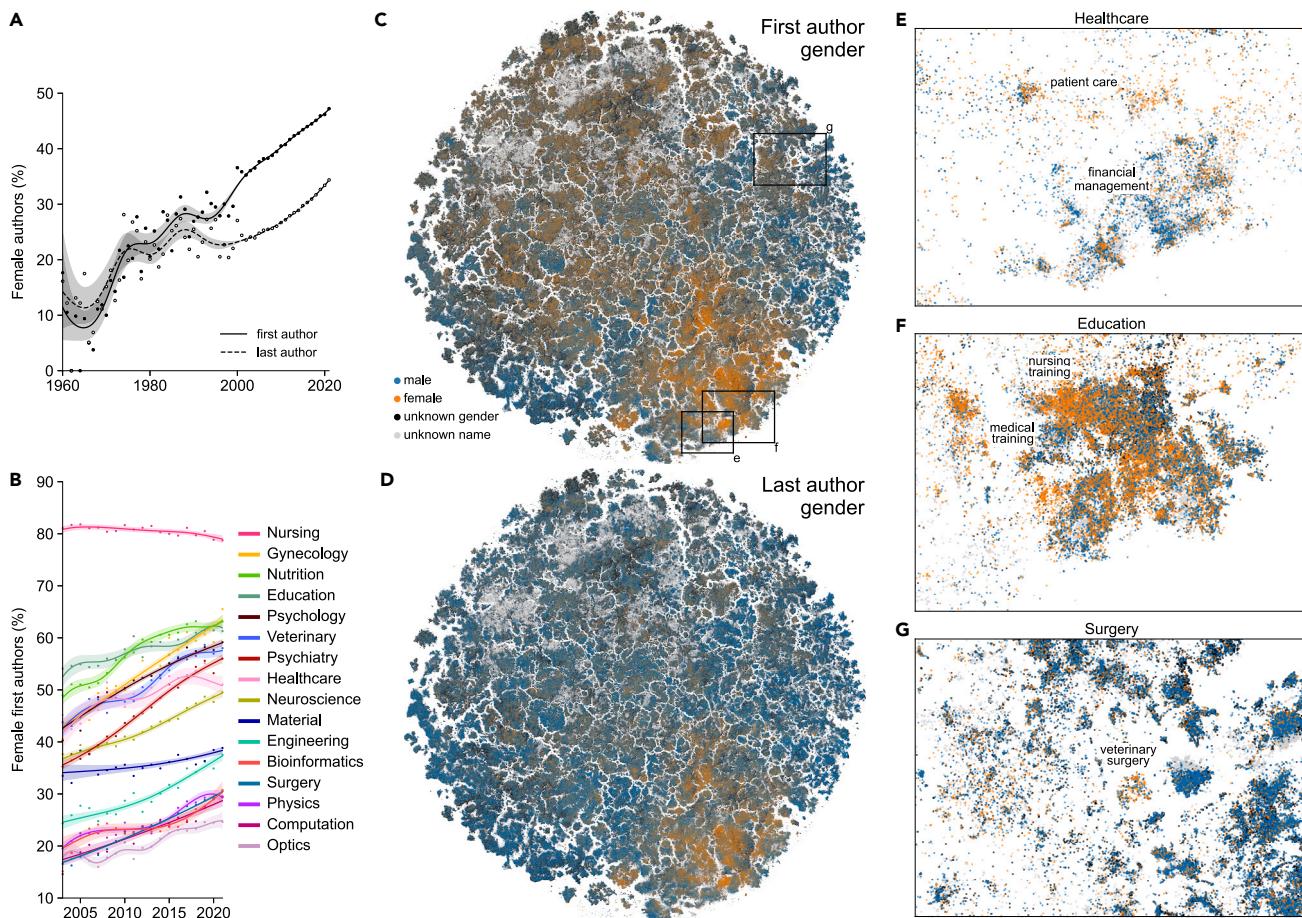


Figure 5. Gender bias in academic authorship

(A) Fraction of female first and last authors across time. The amount of available first names increased dramatically after 2003 (Figure S1C). Smooth curves and confidence intervals were obtained using generalized additive models (see [experimental procedures](#)).

(B) Fraction of female first authors across time for different disciplines.

(C) Papers colored by the inferred gender of their first authors.

(D) Papers colored by the inferred gender of their last authors.

(E-G) Regions of the map showing within-label heterogeneity in the distribution of first authors' gender: in healthcare (E), education (F), and surgery (G). Only papers belonging to those labels are shown.

accuracy, the better the embedding, and found this metric to be useful to compare different representations (Tables 1 and S2).

BERT-based models

We used PubMedBERT³ to obtain a numerical representation of each abstract. Specifically, we used the HuggingFace's transformers library and the publicly released PubMedBERT model. PubMedBERT is a Transformer-based language model trained in 2020 on PubMed abstracts and full-text articles from PubMed Central.

In pilot experiments, we compared performance of eight BERT variants: the original BERT,¹⁰ SciBERT,¹² BioBERT,¹³ PubMedBERT,³ SBERT,¹¹ SPECTER,¹⁴ SimCSE,¹⁵ and SciNCL.¹⁶ The exact HuggingFace models that we used were:

- (1) bert-base-uncased
- (2) allenai/scibert_scivocab_uncased
- (3) dmis-lab/biobert-v1.1
- (4) microsoft/BioMedNLP-PubMedBERT-base-uncased-abstract-fulltext
- (5) sentence-transformers/all-mpnet-base-v2
- (6) allenai/specter

(7) malteos/scinl

(8) princeton-nlp/unsup-simcse-bert-base-uncased

All of these models have the same architecture (bert-base; 110 million parameters) but were trained and/or fine-tuned on different data. The original BERT was trained on a corpus of books and text from Wikipedia. SciBERT was trained on a corpus of scientific articles from different disciplines. BioBERT fine-tuned the original BERT on PubMed abstracts and full-text articles from PubMedCentral. PubMedBERT was trained on the same data from scratch (and its vocabulary was constructed from PubMed data, whereas BioBERT used BERT's vocabulary).

The other four models were fine-tuned to produce sentence embeddings instead of word embeddings, i.e., to generate a single vector representation of the entire input text (we treated each entire abstract as one single "sentence" when providing it to these models). SBERT fine-tuned BERT using a corpus of similar sentences and paragraphs; the specific model that we used was obtained via fine-tuning MPNet.⁵⁴ According to SBERT's authors, this is currently the most powerful generic SBERT model; note that their training procedure has evolved since the original approach described in Reimers and Gurevych.¹¹ SPECTER and SciNCL, both fine-tuned the SciBERT

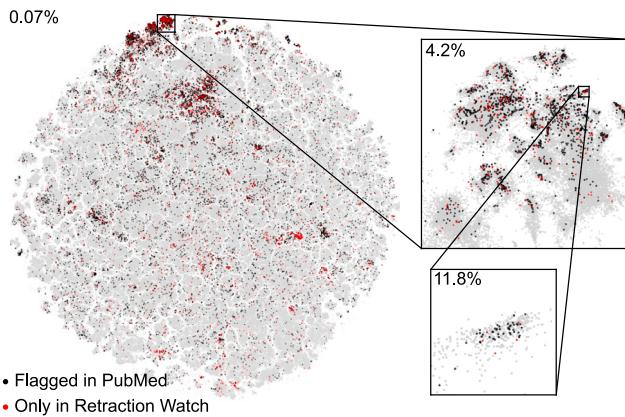


Figure 6. Retracted papers group together

All papers flagged as retracted by PubMed with intact abstracts (11,756) are highlighted in black, plotted on top of the non-retracted papers. Additional retracted papers (3,572) from the Retraction Watch database are shown in red. First inset corresponds to one of the regions with higher density of retracted papers (4.2%), covering research on cancer-related drugs, marker genes, and microRNA. Second inset corresponds to a subregion with a particularly high fraction of retracted papers (11.8%), the one we used for manual inspection.

model using contrastive loss functions based on the citation graph. SimCSE fine-tuned the original BERT using a contrastive loss function between the sentence representations obtained with two different dropout patterns, using Wikipedia texts.

For this pilot experiment, we used a subset of our data ($n = 1,000,000$ labeled papers; 990,000 were used as a training set and 10,000 as a test set) to measure kNN accuracy ($k = 10$) of each of these models, and obtained the highest accuracy with PubMedBERT (see Table S2). This made sense as PubMedBERT's training data largely overlapped with our dataset. We found that SBERT performed better than BERT, but did not reach the level of PubMedBERT on our task. SimCSE did not outperform the original BERT in our benchmark. SPECTER and SciNCL outperformed SciBERT, suggesting that citation information can be helpful for training scientific language models. Still, both models performed worse than PubMedBERT on our task.

Furthermore, we compared kNN accuracy after t-SNE between different BERT models (Figure S12), and again obtained the best results with PubMedBERT (Table S3). The performance of SciNCL here was only 0.1% lower. We used the same settings for t-SNE as described below, but ran it with the default number of iterations (750).

Each abstract gets split into a sequence of tokens, and PubMedBERT represents each token in a 768-dimensional latent space. PubMedBERT's maximum input length is 512 tokens and longer abstracts are automatically truncated at 512 tokens (this corresponds to roughly 300–400 words, and ~98% of all abstracts were shorter than 512 tokens). We are interested in a single 768-dimensional representation of each abstract, rather than 512 of them. For this, we compared several approaches commonly used in the literature: using the representation of the initial [CLS] token, the trailing [SEP] token, and averaging the representations of all tokens.^{10–12} Using the [SEP] token yielded the highest kNN accuracy in our pilot experiments (Table S2), so we adopted this approach.

Note that sentence transformers were originally trained to optimize one specific representation, e.g., SBERT uses the average representation across all tokens as its sentence-level output, while SPECTER uses the [CLS] token. For consistency, in Table S2 we report the performance of all three representations for each model. SBERT implementation (sentence-transformers library) normalizes its output to have norm 1. In Table S2 we report the accuracy without this normalization (64.5%), as obtained using the transformers library; with normalization, the accuracy changed by less than 0.1%.

Su et al.⁵⁵ argued that whitening BERT representation can lead to a strongly improved performance on some benchmarks. We tried whitening the PubMedBERT representation, but only observed a decrease in the kNN accuracy. For this experiment, we used a test set of 500 labeled papers, and compared PubMedBERT without any transformations, after centering, and after whitening, using both Euclidean metric and the cosine metric, following Su et al.⁵⁵ We obtained the best results using the raw PubMedBERT representation (Table S4). Our conclusion is that whitening does not improve the kNN graph of the PubMedBERT representation.

In the end, our entire collection of abstracts is represented as a 20,687,150×768 dense matrix.

TF-IDF representation

In our prior work,¹⁸ we used the bag-of-words representation of PubMed abstracts and compared several different normalization approaches. We obtained the highest kNN accuracy using the TF-IDF representation⁵⁶ with log-scaling, as defined in the scikit-learn implementation (version 0.24.1):

$$X_{ij} = (1 + \ln C_{ij}) \cdot \left(1 + \ln \frac{1+n}{1+\sum_k(C_{kj}>0)} \right)$$

if $C_{ij} > 0$ and $X_{ij} = 0$ otherwise. Here, n is the total number of abstracts and C_{ij} are word counts, i.e., the number of times word j occurs in abstract i . In the scikit-learn implementation, the resulting X_{ij} matrix is then row-normalized, so that each row has ℓ_2 norm equal to 1.

This results in a 20,687,150×4,679,130 sparse matrix with 0.0023% non-zero elements, where 4,679,130 is the total number of unique words in all abstracts.

This matrix is too large to use in t-SNE directly, so for computational convenience we used truncated SVD (sklearn.decomposition.TruncatedSVD with algorithm = “arpack”) to reduce dimensionality to 300, the largest dimensionality we could obtain given our RAM resources. Note that we did not use SVD when using BERT representations and worked directly with 768-dimensional representations.

The kNN accuracy values for the TF-IDF and SVD ($d = 300$) representations measured on the same 1 million subset as used in the previous section were 61.0% and 54.8%, respectively. After t-SNE, the kNN accuracy was 49.9%. After our analysis has already been completed, we tried row-normalizing the SVD representation and observed that this increased the kNN accuracy to 58.7% (and 52.0% after t-SNE); this is equivalent to using cosine distance instead of Euclidean distance for finding nearest neighbors.

We have also experimented with constructing a TF-IDF representation based on the PubMedBERT's tokenizer, instead of the default TF-IDF tokenizer. This reduces the vocabulary size and the dimensionality of the resulting space from 758,111 to 29,047 (because PubMedBERT's tokenizer does not include all unique words as tokens, but instead fragments rare words into repeating substrings). This barely affected kNN classification accuracy: it changed from 61.0% to 61.6% in the high-dimensional space and from 49.9% to 50.1% in the 2D space after SVD and t-SNE. Note that the actual PubMedBERT representation captures many more aspects of the text than just the presence or absence of specific tokens, so it is unsurprising that the representation quality was higher there (67.7% in 768D and 60.8% in 2D).

t-SNE

We used the openTSNE (version 0.6.0) implementation⁵⁶ of t-SNE⁴ to reduce dimensionality from 768 (for the BERT representation) or 300 (for the TF-IDF representation) to $d = 2$. OpenTSNE is a Python reimplementation of the Fit-SNE⁵⁷ algorithm.

We ran t-SNE following the procedure established in our prior work¹⁸: using uniform affinities (on the approximate kNN graph with $k = 10$) instead of perplexity-based affinities, early exaggeration annealing instead of the abrupt switch of the early exaggeration value, and extended optimization for 2,250 iterations instead of the default 750 (250 iterations for the early exaggeration annealing, followed by 2,000 iterations without exaggeration). We did not use any “late” exaggeration after the early exaggeration phase. All other parameters were kept at default values, including PCA initialization and learning rate set to $n/12$, where n is the sample size.

In our previous work we showed that this visualization approach outperformed UMAP (version 0.5.1)¹⁷ on PubMed data in TF-IDF representation in terms of both k NN recall and k NN accuracy.¹⁸ We confirmed that the same was true for the PubMedBERT representation of the 1M subset used in the previous sections (Figure S13): the UMAP embedding was qualitatively similar to the t -SNE embedding with exaggeration $\rho = 4$, and its k NN recall (2.8%) and accuracy (49.6%) were lower than those obtained using t -SNE without exaggeration (12.5% and 60.8%, respectively).

The t -SNE embeddings of a PubMed subset containing 1 million papers (Figures S12 and S13; Table S3) used the default number of iterations (750).

The embeddings based on the TF-IDF and PubMedBERT representation showed similar large-scale organization. As t -SNE loss function is unaffected by rotations and/or sign flips, we flipped the x and/or y coordinates of the TF-IDF t -SNE embedding to match its orientation to the PubMedBERT t -SNE embedding. The same was done for the embeddings shown in Figures S12 and S13.

Performance metrics

All k NN-based metrics were based on $k = 10$ exact nearest neighbors, obtained using the NearestNeighbors and KNeighborsClassifier classes from scikit-learn (version 1.0.2) using algorithm = "brute" and n_jobs = -1.⁵⁸

To predict each test paper's label, the k NN classifier takes the majority label among the paper's nearest neighbors in the training set. To measure the accuracy, the classifier was trained on all labeled papers excluding a random test set of labeled papers. The test set size was 5,000 for the high-dimensional representations and 10,000 for the 2D ones. The chance-level k NN accuracy was obtained using the DummyClassifier from scikit-learn with strategy = "stratified," and test set size 10,000.

To predict each test paper's publication year, we took the average publication year of the paper's nearest neighbors in the training set. To measure the root mean-squared error (RMSE), we used the training set consisting of all papers excluding a random test set. The test set size was 5,000 for the high-dimensional representations and 10,000 for the 2D ones. The chance-level RMSE was calculated by drawing 10 random papers instead of nearest neighbors, for a test set of 5,000 papers.

We define k NN recall as the average size of the overlap between k nearest neighbors in the high-dimensional space and k nearest neighbors in the low-dimensional space. We averaged the size of the overlap across a random set of 10,000 papers for the BERT representation and 5,000 papers for the TF-IDF representation. The k NN recall value reported in Table 1 for the TF-IDF representation measures the recall of the original TF-IDF neighbors (0.7%); the recall of the neighbors from the SVD space (which was used for t -SNE) was 1.5%.

Isolatedness metric was defined as the average fraction of k nearest neighbors belonging to the same corpus. We used a random subset of 5,000 papers from each corpus to estimate the isolatedness. The regions from Table 2 were selected as follows. The HIV/AIDS set contained all papers with "HIV" or "AIDS" words (upper case or lower case) appearing in the abstract. The influenza set contained all papers with the word "influenza" in the abstract (capitalized or not). Similarly, the meta-analysis set was obtained using the word "meta-analysis." The virology and ophthalmology sets correspond to the journal-based labels (see above).

COVID-related papers

We considered a paper COVID-related if it contained at least one of the following terms in its abstract: "covid-19," "COVID-19," "Covid-19," "CoVid-19," "2019-nCoV," "SARS-CoV-2," "coronavirus disease 2019," "Coronavirus disease 2019." Our dataset included 132,802 COVID-related papers.

We selected 27 frequent terms contained in COVID-related paper titles to highlight different subregions of the COVID cluster. The terms were: antibody, anxiety, cancer, children, clinical, epidemic, healthcare, immune, implications, mental, mortality, outbreak, pediatric, pneumonia, population, psychological, respiratory, social, strategies, students, surgery, symptoms, therapy, transmission, treatment, vaccine, and workers. Papers were assigned a keyword if their title contained that term, either capitalized or not. Paper titles containing more than one term were assigned randomly to one of them. This resulted in 35,874 COVID-related papers containing one of those keywords: 27.0% from the total amount of COVID-related papers

and 45.6% of the COVID-related papers from the main COVID cluster in the embedding.

Generalized additive models

We used generalized additive models (GAMs) to obtain smooth trends for several of our analyses across time (Figures 3C, 4C, 5C, and 5D). We used the LinearGAM (GAM with the Gaussian error distribution and the identity link function) and the LogisticGAM (GAM with the binomial error distribution and the logit link function) from the pyGAM Python library (version 0.8.0).⁵⁹ In all cases, we excluded papers published in 2022, since we only had very few of them (as we used the 2021 baseline of the PubMed dataset, see above). Linear GAMs (with n_splines = 6) were used for Figure 3C, and logistic GAMs (with n_splines = 12) were used for Figures 4C, 5C, and 5D. All GAMs had the publication year as the only predictor.

In all cases, we used the gridsearch() function to estimate the optimal smoothing (lambda) parameter using cross-validation. To obtain the smooth curves shown in the plots, we predicted the dependent value on a grid of publication years. The confidence intervals were obtained using the confidence_intervals() function from the same package.

In Figure 3C, the response variable was k NN overlap of a neuroscience paper with the target discipline. For each discipline, the input data were a set of 500 randomly chosen neuroscience papers for each year in 1975–2021. If the total number of neuroscience papers for a given year was less than 500, all of them were taken for the analysis. The k NN overlap values of individual papers were calculated using $k = 10$ nearest neighbors obtained with the NearestNeighbors class.

In Figure 4C, the binary response variable was whether a paper contained "machine learning" in its abstract. For each discipline, the input data were all 2010–2021 papers.

In Figures 5C and 5D, the binary response variable was whether the paper's first or last author was female (as inferred by the gender tool, see below). The input data in all cases were all papers with gender information from 1960 to 2021.

Gender inference

We extracted authors' first names from the XML tag ForeName that should in principle only contain the first name. However, we observed that sometimes it contained the full name. For that reason, we always took the first word of the ForeName tag contents (after replacing hyphens with spaces) as the author's first name. This reduced some combined first names (such as Eva-Maria or Jose Maria) to their initial word (Eva; Jose). In many cases, mostly in older papers, the only available information about the first name was an initial. As it is not possible to infer gender from an initial, we discarded all extracted first names with length 1. In the end we obtained 13,429,169 first names of first authors (64.9% of all papers) and 13,189,271 first names of last authors (63.8%), almost only from 1960 to 2022.

We used the R package gender²⁶ (version 0.6.0) to infer authors' genders. This package uses a historical approach that takes into account how naming practices have changed over time, e.g., Leslie used to be a male name in the early twentieth century but later has been mainly used as a female name. For each first/last author, we provided gender with the name and the publication year, and obtained the inferred gender together with a confidence measure.

The gender package offers inference based on different training databases. We used the 1930–2012 Social Security Administration data from the USA (method = "ssa"). For the papers published before 1930 we fixed the year to 1930 and for the papers published after 2012, we fixed it to 2012. The SSA data do not contain information on names that are not common in the USA, and we only obtained inferred genders for 8,363,116 first authors (62.3% of available first names) and 8,468,165 last authors (63.1% of available last names). Out of all inferred genders, 3,543,592 first authors (42.4%) and 2,464,882 last authors (29.1%) were female.

Importantly, our gender inference is only approximate.²⁶ The inference model has clear limitations, including limited US-based training data and state-imposed binary genders. Moreover, some first names are inherently gender-ambiguous. However, the distribution of inferred genders over biomedical fields and the pattern of changes over the last decades matched what is known about the gender imbalance in academia, suggesting that inferred genders were sufficiently accurate for our purposes.

Retracted papers

We obtained PMIDs of papers classified in PubMed as retracted (13,569) using the PubMed web interface on 19.04.2023. Of those, 11,998 were present in our map (the rest were either filtered out in our pipeline or not included in the 2021 baseline dataset we used). To make sure that retracted papers were not grouping together in the BERT space because their abstract had been modified to indicate a retraction, we excluded from consideration all retracted papers containing the words “retracted,” “retraction,” “withdrawn,” or “withdrawal” in their abstract (242 papers). The remaining retracted papers (11,756) had intact original abstracts and are shown in [Figure 6](#).

There was one small island at the bottom of the map containing retraction notices (they have independent PubMed entries with separate PMIDs) as well as corrigenda and errata, which were not filtered out by our length cutoffs. Many of the 242 retracted papers with post-retraction modified abstracts were also located there.

We obtained the Retraction Watch database through (<https://api.labs.crossref.org/data/retractionwatch?name@email.org>) as a CSV file (41 MB) on 21.09.2023. It contained 18,786 retracted papers indexed in PubMed. Of those, 15,666 were present in our map (the rest were either filtered out in our pipeline or not included in the 2021 baseline dataset we used). 15,103 of those were intact papers. These 15,103 papers contained all of the 11,998 papers used above except for 234 papers. This gave 3,572 additional retracted papers shown in [Figure 6](#) in red.

2023 annual PubMed baseline

While our paper was in revision, we updated the dataset by downloading the latest annual PubMed snapshot (2023 baseline; files called pubmed24n0001.xml.gz to 1219.xml.gz, download date: 06.02.2024, 350 GB). We used this entire dataset, and not only the files containing 2022–2023 papers, to avoid duplicated entries and to use the latest metadata. This snapshot included in total 36,555,430 papers. After filtering with our previous criteria, we were left with 23,389,083 papers.

We extracted the same attributes from the metadata as described above, with the addition of the first affiliation of the first author. We used this affiliation to assign each paper to a country ([Figure S10](#)), by searching the string for all existing country names in English (taking into account possible name variations, such as “United Kingdom” and “UK”). Consequently, papers that may have included country names in their original language (e.g., “Deutschland” instead of “Germany”) were not matched to any country. We noticed that many US affiliations did not explicitly include country name so we assigned affiliations containing a name of any US state to the US. This resulted in 19,937,913 papers (85.2%) with an assigned country. We matched the affiliation countries to our main dataset (2021 baseline) using PMIDs, which led to 17,404,977 papers (84.1%) with an assigned country ([Figure S8](#)).

We used the same journal-based labels to color the embedding ([Figure S9](#)) and added “dentistry” as an additional label. This resulted in 8,028,583 labeled papers (34.3%).

Runtimes

Computations were performed on a machine with 384 GB of RAM and Intel Xeon Gold 6226R processor (16 multi-threaded 2.90 GHz CPU cores) and on a machine with 512 GB of RAM and Intel Xeon E5-2630 version 4 processor (10 multi-threaded 2.20 GHz CPU cores). BERT embeddings were calculated using an NVIDIA TITAN Xp GPU with 12.8 GB of RAM.

Parsing the XML files took 10 h, computing the PubMedBERT embeddings took 74 h, running t-SNE took 8 h. More details are given in [Table S5](#). We used exact nearest neighbors for all kNN-based quality metrics, so evaluation of the metrics took longer than computing the embedding. In total, it took around 8 days to compute all the reported metrics ([Table S5](#)).

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.patter.2024.100968>.

ACKNOWLEDGMENTS

We thank Richard Van Noorden, David Bimler, Ivan Oransky, and Jennifer Byrne for discussions. This research was funded by the Deutsche Forschungsgemeinschaft (KO6282/2-1, BE5601/8-1, EXC 2064 “Machine Learning: New Perspectives for Science” 390727645, and EXC 2181 “STRUCTURES” 390900948), by the German Ministry of Education and Research (Tübingen AI Center), and by the Hertie Foundation. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting R.G.-M.

AUTHOR CONTRIBUTIONS

R.G.-M. and D.K. designed the study. R.G.-M. performed the analysis and prepared the figures. L.S. did pilot experiments with language models. B.M.S. developed the interactive website. R.G.-M. wrote the initial draft of the manuscript. R.G.-M., P.B., and D.K. discussed the results and edited the paper. D.K. and P.B. supervised the study.

DECLARATION OF INTERESTS

B.M.S. is Vice President of Information at Nomic AI.

Received: October 31, 2023

Revised: January 16, 2024

Accepted: March 15, 2024

Published: April 9, 2024

REFERENCES

1. Larsen, P.O., and von Ins, M. (2010). The rate of growth in scientific publication and the decline in coverage provided by science citation index. *Scientometrics* 84, 575–603.
2. Bornmann, L., and Mutz, R. (2015). Growth rates of modern science: A bibliometric analysis based on the number of publications and cited references. *J. Assoc. Inf. Sci. Technol.* 66, 2215–2222.
3. Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., and Poon, H. (2021). Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthc.* 3, 1–23.
4. van der Maaten, L., and Hinton, G. (2008). Visualizing data using t-SNE. *J. Mach. Learn. Res.* 9.
5. Boyack, K.W., Smith, C., and Klavans, R. (2020). A detailed open access model of the PubMed literature. *Sci. Data* 7, 408.
6. Börner, K., Klavans, R., Patek, M., Zoss, A.M., Biberstine, J.R., Light, R.P., Larivière, V., and Boyack, K.W. (2012). Design and update of a classification system: The UCSD map of science. *PLoS One* 7, e39464.
7. Nomic, A.I. (2022). Deepscatter. URL: <https://github.com/nomic-ai/deepscatter>
8. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Adv. Neural Inf. Process. Syst.* 30, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., pp. 5998–6008.
9. Sparck Jones, K. (1972). A statistical interpretation of term specificity and its application in retrieval. *J. Doc.* 28, 11–21.
10. Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–4186.
11. Reimers, N., and Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992.
12. Beltagy, I., Lo, K., and Cohan, A. (2019). SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on*

- Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3615–3620.
13. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., and Kang, J. (2020). BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 1234–1240.
 14. Cohan, A., Feldman, S., Beltagy, I., Downey, D., and Weld, D.S. (2020). Specter: Document-level representation learning using citation-informed transformers. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 2270–2282.
 15. Gao, T., Yao, X., and Chen, D. (2021). SimCSE: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 6894–6910.
 16. Ostendorff, M., Rethmeier, N., Augenstein, I., Gipp, B., and Rehm, G. (2022). Neighborhood contrastive learning for scientific document representations with citation embeddings. In The 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP 2022) (Association for Computational Linguistics), pp. 11670–11688.
 17. McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform Manifold Approximation and Projection for dimension reduction. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1802.03426>.
 18. González-Márquez, R., Berens, P., and Kobak, D. (2022). Two-dimensional visualization of large document libraries using t-SNE. In *Proceedings of Topological, Algebraic, and Geometric Learning Workshops 2022* vol. 196 of *Proceedings of Machine Learning Research* (PMLR), pp. 133–141.
 19. Ke, Q. (2019). Identifying translational science through embeddings of controlled vocabularies. *J. Am. Med. Inf. Assoc.* 26, 516–523.
 20. Topol, E.J. (2019). High-performance medicine: the convergence of human and artificial intelligence. *Nat. Med.* 25, 44–56.
 21. Filardo, G., Da Graca, B., Sass, D.M., Pollock, B.D., Smith, E.B., and Martinez, M.A.-M. (2016). Trends and comparison of female first authorship in high impact medical journals: observational study (1994–2014). *BMJ* 352, i847.
 22. Larivière, V., Ni, C., Gingras, Y., Cronin, B., and Sugimoto, C.R. (2013). Bibliometrics: Global gender disparities in science. *Nature* 504, 211–213.
 23. Shen, Y.A., Webster, J.M., Shoda, Y., and Fine, I. (2018). Persistent underrepresentation of women's science in high profile journals. Preprint at bioRxiv. <https://doi.org/10.1101/275362>.
 24. Dworkin, J.D., Linn, K.A., Teich, E.G., Zurn, P., Shinohara, R.T., and Bassett, D.S. (2020). The extent and drivers of gender imbalance in neuroscience reference lists. *Nat. Neurosci.* 23, 918–926.
 25. Bendels, M.H.K., Müller, R., Brueggemann, D., and Groneberg, D.A. (2018). Gender disparities in high-quality research revealed by Nature Index journals. *PLoS One* 13, e0189136.
 26. Blevins, C., and Mullen, L. (2015). Jane, John... Leslie? A historical method for algorithmic gender prediction. *DHQ: Digital Humanities Quarterly* 9, 000223.
 27. Byrne, J.A., and Labbé, C. (2017). Striking similarities between publications from China describing single gene knockdown experiments in human cancer cell lines. *Scientometrics* 110, 1471–1493.
 28. Byrne, J.A., Grima, N., Capes-Davis, A., and Labbé, C. (2019). The possibility of systematic research fraud targeting under-studied human genes: causes, consequences, and potential solutions. *Biomark. Insights* 14, 1177271919829162.
 29. Candal-Pedreira, C., Ross, J.S., Ruano-Ravina, A., Egilman, D.S., Fernández, E., and Pérez-Ríos, M. (2022). Retracted papers originating from paper mills: cross sectional study. *BMJ* 379, e071517.
 30. Bielack, S.S., and Palmerini, E. (2022). A special jubilee: 100 fake osteosarcoma articles. *ESMO open* 7, 100358.
 31. Byrne, J. (2019). We need to talk about systematic fraud. *Nature* 566, 9–10.
 32. Byrne, J.A., and Christopher, J. (2020). Digital magic, or the dark arts of the 21st century — how can journals and peer reviewers detect manuscripts and publications from paper mills? *FEBS Lett.* 594, 583–589.
 33. Else, H., and Van Noorden, R. (2021). The fight against fake-paper factories that churn out sham science. *Nature* 591, 516–519.
 34. Zhao, T., Dai, T., Lun, Z., and Gao, Y. (2021). An analysis of recently retracted articles by authors affiliated with hospitals in mainland China. *J. Sch. Publish.* 52, 107–122.
 35. Fanelli, D., Schleicher, M., Fang, F.C., Casadevall, A., and Bik, E.M. (2022). Do individual and institutional predictors of misconduct vary by country? Results of a matched-control analysis of problematic image duplications. *PLoS One* 17, e0255334.
 36. Sabel, B.A., Knaack, E., Gigerenzer, G., and Bilc, M. (2023). Fake publications in biomedical science: Red-flagging method indicates mass production. Preprint at medRxiv. <https://doi.org/10.1101/2023.05.06.23289563>.
 37. Oransky, I., Fremes, S.E., Kurlansky, P., and Gaudino, M. (2021). Retractions in medicine: the tip of the iceberg. *Eur. Heart J.* 42, 4205–4206.
 38. Wattenberg, M., Viégas, F., and Johnson, I. (2016). How to use t-SNE effectively. *Distill* 1, e2.
 39. Kobak, D., and Berens, P. (2019). The art of using t-SNE for single-cell transcriptomics. *Nat. Commun.* 10, 5416.
 40. Böhm, J.N., Berens, P., and Kobak, D. (2022). Attraction-repulsion spectrum in neighbor embeddings. *J. Mach. Learn. Res.* 23, 1–32.
 41. Kobak, D., and Linderman, G.C. (2021). Initialization is critical for preserving global data structure in both t-SNE and UMAP. *Nat. Biotechnol.* 39, 156–157.
 42. Xu, J., Kim, S., Song, M., Jeong, M., Kim, D., Kang, J., Rousseau, J.F., Li, X., Xu, W., Torvik, V.I., et al. (2020). Building a PubMed knowledge graph. *Sci. Data* 7, 205.
 43. Schmidt, B. (2018). Stable random projection: Lightweight, general-purpose dimensionality reduction for digitized libraries. *Journal of Cultural Analytics* 3, 11033.
 44. Grover, A., and Leskovec, J. (2016). node2vec: Scalable feature learning for networks. In Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 855–864.
 45. Noichl, M. (2021). Modeling the structure of recent philosophy. *Synthese* 198, 5089–5100.
 46. Priem, J., Piwowar, H., and Orr, R. (2022). OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2205.01833>.
 47. Sinha, A., Shen, Z., Song, Y., Ma, H., Eide, D., Hsu, B.-J., and Wang, K. (2015). An Overview of Microsoft Academic Service (MAS) and Applications. In Proceedings of the 24th international conference on world wide web, pp. 243–246.
 48. Liu, F., Vulić, I., Korhonen, A., and Collier, N. (2021). Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 1442–1459.
 49. Yasunaga, M., Leskovec, J., and Liang, P. (2022). LinkBERT: Pretraining language models with document links. In Association for Computational Linguistics (ACL).
 50. Luo, R., Sun, L., Xia, Y., Qin, T., Zhang, S., Poon, H., and Liu, T.-Y. (2022). BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings Bioinf.* 23, bbac409.
 51. Stanford, C.R.F.M. (2022). MosaicML. BioMedLM. URL: <https://huggingface.co/stanford-crfm/BioMedLM>
 52. González-Márquez, R., Schmidt, L., Schmidt, B.M., and Berens, P. (2024). Kobak, D. berenslab/pubmed-landscape: Submission v1. <https://doi.org/10.5281/zenodo.1072758>.
 53. González-Márquez, R., Schmidt, L., Schmidt, B.M., Berens, P., and Kobak, D. (2023). Data from the paper "The landscape of biomedical research". <https://doi.org/10.5281/zenodo.7849020>.

54. Song, K., Tan, X., Qin, T., Lu, J., and Liu, T.-Y. (2020). Mpnet: Masked and permuted pre-training for language understanding. *Adv. Neural Inf. Process. Syst.* **33**, 16857–16867.
55. Su, J., Cao, J., Liu, W., and Ou, Y. (2021). Whitening sentence representations for better semantics and faster retrieval. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2103.15316>.
56. Poličar, P.G., Stražar, M., and Zupan, B. (2019). openTSNE: a modular Python library for t-SNE dimensionality reduction and embedding. Preprint at bioRxiv. <https://doi.org/10.1101/731877>.
57. Linderman, G.C., Rachh, M., Hoskins, J.G., Steinerberger, S., and Kluger, Y. (2019). Fast interpolation-based t-SNE for improved visualization of single-cell RNA-seq data. *Nat. Methods* **16**, 243–245.
58. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830.
59. Servén, D., and Brummitt, C. (2018). pyGAM: Generalized additive models in python. <https://doi.org/10.5281/zenodo.1208723>.