

# Linguistic Analysis of the bioRxiv Preprint Landscape

This manuscript ([permalink](#)) was automatically generated from [greenelab/annoxiver manuscript@84a2c3c](#) on January 29, 2021.

## Authors

---

- **David N. Nicholson**

 [0000-0003-0002-5761](#) ·  [danich1](#)

Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine University of Pennsylvania, Philadelphia PA, USA · Funded by The Gordon and Betty Moore Foundation (GBMF4552); The National Institutes of Health (T32 HG000046)

- **Vincent Rubinetti**

·  [vincerubinetti](#) ·  [vincerubinetti](#)

Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine University of Pennsylvania, Philadelphia PA, USA · Funded by The Gordon and Betty Moore Foundation (GBMF4552); The National Institutes of Health (T32 HG010067)

- **Dongbo Hu**

·  [dongbohu](#)

Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine University of Pennsylvania, Philadelphia PA, USA · Funded by The Gordon and Betty Moore Foundation (GBMF4552); The National Institutes of Health (T32 HG010067)

- **Marvin Thielk**

 [0000-0002-0751-3664](#) ·  [MarvinT](#) ·  [TheNeuralCoder](#)

Elsevier

- **Lawrence E. Hunter**

 [0000-0003-1455-3370](#) ·  [LEHunter](#) ·  [ProfLHunter](#)

Center for Computational Pharmacology, University of Colorado School of Medicine, Aurora CO, USA · Funded by The Gordon and Betty Moore Foundation (GBMF4552)

- **Casey S. Greene**

 [0000-0001-8713-9213](#) ·  [cgreene](#) ·  [greenescientist](#)

Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine University of Pennsylvania, Philadelphia PA, USA; Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Aurora CO, USA · Funded by The Gordon and Betty Moore Foundation (GBMF4552); The National Institutes of Health (T32 HG010067)

# Abstract

---

Preprints allow researchers to make their findings available to the scientific community before they have undergone peer review. This provides an opportunity to understand life sciences peer review and publishing practices. Studies of bioRxiv to date have been largely focused on article metadata, and have examined how preprints are downloaded, cited, published, and discussed online. A missing element has been examining the language contained within preprints. Comparing preprints with peer reviewed manuscripts provides an opportunity to examine how peer review changes these documents.

We sought to compare and contrast linguistic features within bioRxiv preprints to biomedical text as a whole and also their published counterparts. The most prevalent elements that changed appeared to be associated with typesetting and mentions of supplementary or additional files. We created document embeddings derived from a bioRxiv-trained word2vec model and found that they revealed different scientific approaches, linked unannotated preprint-peer reviewed article pairs, and were associated with the journal where a preprint would eventually be published. We examined factors associated with the time elapsed between the posting of a first preprint and the appearance of a peer reviewed publication and found that preprints with more versions posted and also those with more textual changes, assessed via document embeddings, took longer to publish. The full text landscape of bioRxiv provides an opportunity to explore scientific language and how peer review alters articles. We provide a web application with which users can explore how a bioRxiv or medRxiv preprint is positioned within the landscape and identify the journals and articles that are most similar.

## Introduction

---

The dissemination of research findings is key to science, and initially much of this communication happened orally [1]. During the 17th century, the predominate form of communication shifted to personal letters that were shared from one scientist to another [1]. Scientific journals didn't become a predominant mode of communication until the 19th and 20th centuries when the first journal abstract was created [???,1,2]. Although scientific journals became the primary method of communication, they added high maintenance costs and long publication times to scientific discourse [???,2]. Scientists' solution to some of these issues was to communicate through preprints, which are scholarly works that have yet to undergo peer review process [3,4].

Preprints are commonly hosted on online repositories, where users have open and easy access to these works. Notable repositories include arXiv [5], bioRxiv [6] and medRxiv [7]; however, there are over 60 different repositories available [8]. The burgeoning uptake of preprints in life sciences has been examined through research focused on metadata from the bioRxiv repository. For example, scientists found that life science preprints are being posted at an increasingly high rate [9]. Furthermore, these preprints are being rapidly shared on social media, routinely downloaded, and cited [10]. Certain preprint categories are read and shared by both scientists and non-scientists alike [11]. Overall, about two-thirds to three-quarters of preprints are eventually published [12,13] and life science articles that have a corresponding preprint version are cited and discussed more often than articles without them [14,15,16]. Preprints take an average of 160 days to become published [17], and those with multiple versions take longer to publish [17].

In spite of the success and excitement of preprints, there are number of issues that arose from their constant use. Preprint repositories receive a growing number of submissions, which pose challenges for this mode of communication [18]. For example, repositories have a hard time searching and linking preprints with their published counterparts [15,19], which results in missing links and consequently erroneous metadata. Furthermore, these repositories lack tools to show how textual content of preprints are altered due to the peer review process [18]. These repositories are open

access for all to view preprints, which results in concern from scientists that they could be scooped by competitors [18,20]. Plus, preprint repositories do not have in-depth peer review which can result in posted preprints containing inconsistent results or conclusions [16,19,21,22]. Despite a growing emphasis on using the study of preprints to examine the publishing process in the life sciences, how these findings related to the text of documents within bioRxiv has not been examined.

Textual analysis is a methodology that uses linguistic, statistical and machine learning techniques to analyze and extract information from text [23]. This set of techniques have made a sizable impact within the life science community by providing valuable insight on biomedical text. For instance, scientists analyzed linguistic similarities and differences of biomedical corpora [24,25]. Scientists have provided the community with a number of tools that aide future text mining systems [26,27,28] as well as advice on how to train and test future text processing systems [29,30,31]. We use textual analysis to examine the bioRxiv repository, placing a particular emphasis on understanding the extent to which full text analysis can address hypotheses derived from the analysis of metadata alone.

Preprints are still an emerging method of scientific communication in the life sciences. To understand how they relate to the traditional publishing ecosystem, we examine the linguistic similarities and differences between preprints and peer reviewed text and observe how linguistic features change during the peer review and publishing process. We theorize that preprints and biomedical text would be quite similar, especially when controlling for the differential uptake of preprints across fields. In other contexts, neural networks trained in certain ways can produce a representation of words and documents that has useful properties for many tasks - termed word or document embeddings [32,33]. Here, we hypothesize that using these networks to embedded preprints provides a versatile way to disentangle linguistic features and serves as a practical medium for improving preprint repository functionality. We test this hypothesis by producing a linguistic landscape of bioRxiv preprints, detecting preprints that change substantially during publication, and identify journals that publish manuscripts that are linguistically similar to a target preprint. We encapsulate our findings through a web-app that projects a user-selected preprint onto this landscape and suggests journals and articles that are linguistically similar. Taken together, our work reveals how linguistically similar and dissimilar preprints are to peer reviewed text and quantifies linguistic changes that occur during the peer review process. Furthermore, our work examines the association between linguistic changes and a preprint's time to publication suggesting that comparisons between embeddings of preprints and peer reviewed documents provides a means to study the process of peer review.

## Materials and Methods

---

### Corpora Examined

#### BioRxiv Corpus

BioRxiv [6] is a repository for life sciences preprints. We downloaded an xml snapshot of this repository on February 3, 2020 from bioRxiv's Amazon S3 bucket [34]. This snapshot contained the full text and image content of 98,023 preprints. Preprints on bioRxiv are versioned, and in our snapshot 26,905 out of 98,023 contained more than one version. When preprints had multiple versions, we used the latest one unless otherwise noted. Authors submitting preprints to bioRxiv select one of twenty-nine different categories. Researchers also select an article type, which can be a new result, confirmatory finding, or contradictory finding. Some preprints in this snapshot were withdrawn from bioRxiv: when this happens, their content is replaced with the reason for withdrawal. As there were very few withdrawn preprints, we did not treat these as a special case.

#### PubMed Central Open Access Corpus

PubMed Central (PMC) is a digital archive for the United States National Institute of Health's Library of Medicine (NIH/NLM) that contains full text biomedical and life science articles [35]. PMC articles can be closed access ones from research funded by the NIH appearing after an embargo period or those published under Gold Open Access [36] publishing schemes. Paper availability within PMC is largely dependent on the journal's participation level [37]. Individual journals can fully participate in submitting articles to PMC, selectively participate sending only a few papers to PMC, only submit papers according to NIH's public access policy [38], or not participate at all. As of September 2019, PMC had 5,725,819 articles available [39]. Out of these 5 million articles, about 3 million were open access (PMCOA) and available for text processing systems [27,40]. PMC also contains a resource that holds author manuscripts that have already passed the peer review process [41]. Since these manuscripts have already been peer reviewed, we kept them out of our analysis as the scope of our work is solely focused on examining the beginning and endpoints of a preprint's lifecycle. We downloaded a snapshot of the PMCOA corpus on January 31, 2020. This snapshot contained many types of papers: literature reviews, book reviews, editorials, case reports, research articles and more. We used only research articles, which aligns with the intended role of bioRxiv, and we refer to these articles as the PMCOA Corpus.

## The New York Times Annotated Corpus

The New York Times Annotated Corpus (NYTAC) is [42] is collection of newspaper articles from the New York Times dating from January 1, 1987 to June 19, 2007. This collection contains over 1.8 million articles where 1.5 million of those articles have undergone manual entity tagged by library scientists [42]. We downloaded this collection on August 3rd, 2020 from the Linguistic Data Consortium (see Software and Data Availability section) and used the entire collection as a negative control for our corpora comparison analysis.

## Mapping bioRxiv preprints to their published counterparts

We used CrossRef [43] to identify bioRxiv preprints that were linked to a corresponding published article. We accessed CrossRef on July 7th, 2020 and were able to successfully link 23,271 preprints to their published counterparts. Out of those 23,271 preprint-published pairs only 17,952 pairs had a published version present within the PMCOA corpus. For our analyses that involved published links we only focused on the subset of preprints-published pairs that contained a published article within PMCOA.

## Comparing Corpora

We compared the bioRxiv, PMC, and NYTAC corpora to assess the similarities and differences between them. We use the NYTAC as an out-group to assess the similarity of two life sciences repositories when compared with non-life sciences text. The corpora contain both words and non-word symbols (e.g.,  $\pm$ ), which we refer to together as tokens to avoid confusion. We calculated the following statistics for each corpus: the number of documents, the number of sentences, the total number of tokens, the number of stopwords, the average length of a document, the average length of a sentence, the number of negations, the number of coordinating conjunctions, the number of pronouns and the number of past tense verbs. Next, we used spaCy's "en\_core\_web\_sm" model [44] (version 2.2.3) to preprocess all corpora and filtered out 326 spaCy-provided stopwords.

Following cleaning, we calculated the frequency of every token across all corpora. Because many tokens were unique to one set or the other and observed at low frequency, we used the union of the top 100 most frequent tokens from each pair of corpora to compare them. We generated a contingency table for each token in this union and calculated the odds ratio and 95% confidence

interval [45]. We measured corpus similarity by calculating the Kullback–Leibler divergence across all three corpora, which focuses on token distribution differences as opposed to token-level differences.

## Constructing a Document Representation for Life Sciences Text

We sought to build a model that would capture the linguistic similarity of articles. Word2vec is a suite of neural networks designed to model linguistic features of words based on their appearance in text. These models are trained to either predict a word based on its sentence context as a continuous bag of words (CBOW) or predict the context based on a given word in a skipgram model [32]. Through these prediction tasks the networks learn latent features that can be used for downstream tasks such as identifying similar words. We used gensim [46] (version 3.8.1) to train a word2vec continuous bag of words (CBOW) [32] model over the bioRxiv corpus. Our neural network architecture had 300 hidden nodes, and we trained this model for 20 epochs. We set a fixed random seed and used gensim's default settings for all other hyperparameters. Following training, we generated a document vector for every article within bioRxiv and PubMed Central. We calculated the document vector by taking the average of every token present within a given article [33]. Words absent from the word2vec model were ignored.

## Visualizing and Characterizing Preprint Representations

We sought to visualize the landscape of preprints and determine the extent to which their representation as document vectors corresponded to author-supplied document labels. We used principal component analysis (PCA) [47] to project bioRxiv document vectors into a low dimensional space. We trained this model using the scikit-learn [48] implementation of a randomized solver [49] with a random seed of 100, output of 50 principal components (PCs), and default settings for all other hyperparameters. After fitting, each preprint has a score for each PC. To visualize the tokens associated with each PC, we calculated the cosine similarity of each PC to all tokens in our word2vec model's vocabulary. We report the top 100 positive and negative scoring tokens in the form of word clouds, where the size of each word corresponds to the magnitude of similarity and color represents positive (orange) or negative (blue) association.

## Discovering Unannotated Preprint-Publication Relationships

The bioRxiv maintainers have automated procedures to link preprints to peer reviewed versions and many journals require authors to update preprints with a link to the published version. However, this automation is largely based on exact matching of certain attributes, and authors can forget to establish a link after publication. Authors can change the title between a preprint and published version (e.g., [50] and [51]), which prevents bioRxiv from automatically establishing a link. If the authors do not report the publication to bioRxiv, the preprint and the published version are treated as distinct entities despite representing the same underlying research. We recognized that close proximity in the embedding space could reveal preprint to published version links that were missed by existing automated processes. We used the subset of paper-preprint pairs annotated in CrossRef as described above to calculate the distribution of known preprint to published distances, which we calculated as the Euclidean distance between the preprint's embedding coordinates and the coordinates of its corresponding published version. We also calculated a background distribution, which consisted of the distance between each preprint with an annotated publication and a randomly selected article from the same journal. Next, we calculated distances between preprints without a published version link with PubMed Central articles that weren't matched with a corresponding preprint. We filtered any potential links with distances that were greater than the minimum value of the background distribution to reduce the curation burden. Lastly, we binned the remaining pairs based on percentiles from the annotated pairs distribution at the [0,25th percentile), [25th percentile, 50th percentile), [50th percentile, 75th percentile), and [75th percentile, minimum background



distance). We randomly sampled 50 articles from each bin for manual annotation. We shuffled these four sets to produce a list of 200 potential preprint-published pairs with a randomized order. We supplied these pairs to two co-authors to manually determine if each link between a preprint and a putative matched version was correct or incorrect. After the curation process, we encountered eight disagreements between the reviewers. We supplied the preprint-publication pairs on which reviewers disagreed to a third scientist, who carefully reviewed each case and made a final determination. We used this curated set to evaluate the extent to which distance in the embedding space revealed true but unannotated links between preprints and their published versions.

## Measuring Time Duration for Preprint Publication Process

We measured the time required for preprints to be published in the peer reviewed literature and compared this time within fields and as a function of the extent to which documents changed between the preprint and publication. We queried bioRxiv's application programming interface (API) to obtain the date a preprint was posted onto bioRxiv as well as the date a preprint was accepted for publication. We calculated the difference between the date at which a preprint was first posted and its publication date to provide a publication interval, and we also recorded the number of preprint versions posted onto bioRxiv. To measure the amount of textual difference, we calculated the Euclidean distance between the document representation of each preprint and the corresponding published version. We performed linear regression to model the relationship between preprint version count and a preprint's time to publication as well as the relationship between document representation distances and a preprint's time to publication. We visualized results as square bin plots. We observed a limited number of cases in which authors appeared to post preprints after the date of publication, which results in preprints receiving a negative time difference, as previously reported [52]. We did not remove preprints that had a negative time publication in our linear regression analysis as it was not strictly necessary, but we removed them in our survival curve analysis where they were incompatible with the analytical approach. In practice, the number with negative publication times and the short lead time between publication and preprint has a minimal impact on results.

Document distances can be difficult to understand, so we sought to contextualize the meaning of a distance unit. We selected preprints within the Bioinformatics topic area, which was well-represented on bioRxiv. For preprints submitted to the Bioinformatics topic area, we sampled a pair of preprints and calculated their differences 1000 times and reported the mean.

In addition to contextualizing the document distance, we also wanted to contextualize differences in the time to publication. We examined time to publication for each topic area using the Kaplan-Meier estimator [53] on preprints within bioRxiv, treating preprints not yet published as survival. We generated these curves using the KaplanMeierFitter function from the lifelines [54] (version 0.25.6) python package. We reported the half-life of each bioRxiv preprint category.

## Building Journal Venue Classifiers

We hypothesized that preprints would be more likely to be published in journals that contained similar content to the work in question. To test this hypothesis, we designed an experiment examining document and journal representations. First, we removed all journals that had fewer than 100 papers in the PMC corpus. Certain manuscripts in the PMC corpus were annotated to their corresponding bioRxiv preprints through CrossRef as previously noted. We held out this subset and treated it as a gold standard test set. We used the remainder of the PMC corpus for training and initial evaluation via cross validation using the scikit-learn k-Nearest Neighbors implementation [55]. We imagined a use case of prioritizing relevant journals for preprint authors, and considered a list of ten journal

suggestions to be an appropriate number and we considered a prediction to be a true positive if the correct journal appeared within the ten closest neighbors of the query article.

Certain journals publish articles in a focused topic area, while others publish articles that cover many topics. Likewise, some journals have a publication rate of at most hundreds of papers per year while others publish at a rate of at least ten-thousand papers per year. Accounting for these characteristics, we designed two approaches - one centered on manuscripts and another centered on journals.

For the manuscript-based approach, we identified the ten most similar published manuscripts and evaluated where the documents were published. We embedded each query article into the space defined by the word2vec model as described for preprints. We selected the ten manuscripts that were nearest by Euclidean distance in the embedding space and returned the journal in which they were published. The number of journals returned via this method could be less than ten as multiple papers in close proximity to query article may belong to the same journal. Because this approach was based on paper proximity, we could return the articles that led to each journal being returned. However, journals that publish more papers are more frequently recommended in this framing.

For the journal-based approach, we identified the ten most similar journals by constructing a journal representation in the same embedding space. We computed journal centroids as the average embedding of all published papers in the journal. We then projected a query article into the same space and returned the ten closest journal centroids by Euclidean distance. This technique guaranteed that at least ten distinct journals were returned and prevented journals that publish many papers from being heavily overrepresented.

In both cases, we set the number of neighbors for each model to be 10 and then evaluated both models via 10-fold cross validation. We evaluated performance of both classifiers on our gold standard test set of published preprints.

## Web Application for Discovering Similar Preprints and Journals

We developed a web application that identifies similar papers and journals for any bioRxiv and medRxiv preprint and that places the preprint into the overall document landscape. Our web application downloads a pdf version of a preprint hosted on the bioRxiv or medRxiv server. We use pdfminer [56] to extract text from the downloaded pdf. The extracted text is then fed into our word2vec model to construct a document embedding representation. We pass this representation onto our journal and manuscript search to identify journals based on the ten closest neighbors of individual papers as well as journal centroids. We implemented this search using the scikit-learn implementation of k-d trees. To run it more cost effectively on cloud computing environment, we sharded the k-d trees into four trees.

Accompanying these recommendations, we also provided a neural network derived visualization of our training set and the article's position within it. We used SAUCIE [57], an autoencoder designed to cluster single cell RNA-seq data, to build a two-dimensional embedding space that could be applied to newly generated preprints without retraining, a limitation of other approaches that we explored for visualizing entities expected to lie on a nonlinear manifold. We trained this model on document embeddings of PMC articles that did not contain a matching preprint version. We used the following parameters to train the model: a hidden size of 2, a learning rate of 0.001, lambda\_b of 0, lambda\_c of 0.001, and lambda\_d of 0.001 for 2000 iterations. When a user requests a new document, we can then project the document on the pretrained model to generate a visualization in two-dimensional space. We illustrate our recommendations as a short list and provide access to our network visualization at <https://greenelab.github.io/preprint-similarity-search/>.

We used the fully trained model to project user-requested bioRxiv preprints onto the generated landscape to enable users to see where their preprint falls along the landscape.

## Results

### Comparing bioRxiv to other corpora

#### bioRxiv Metadata Statistics

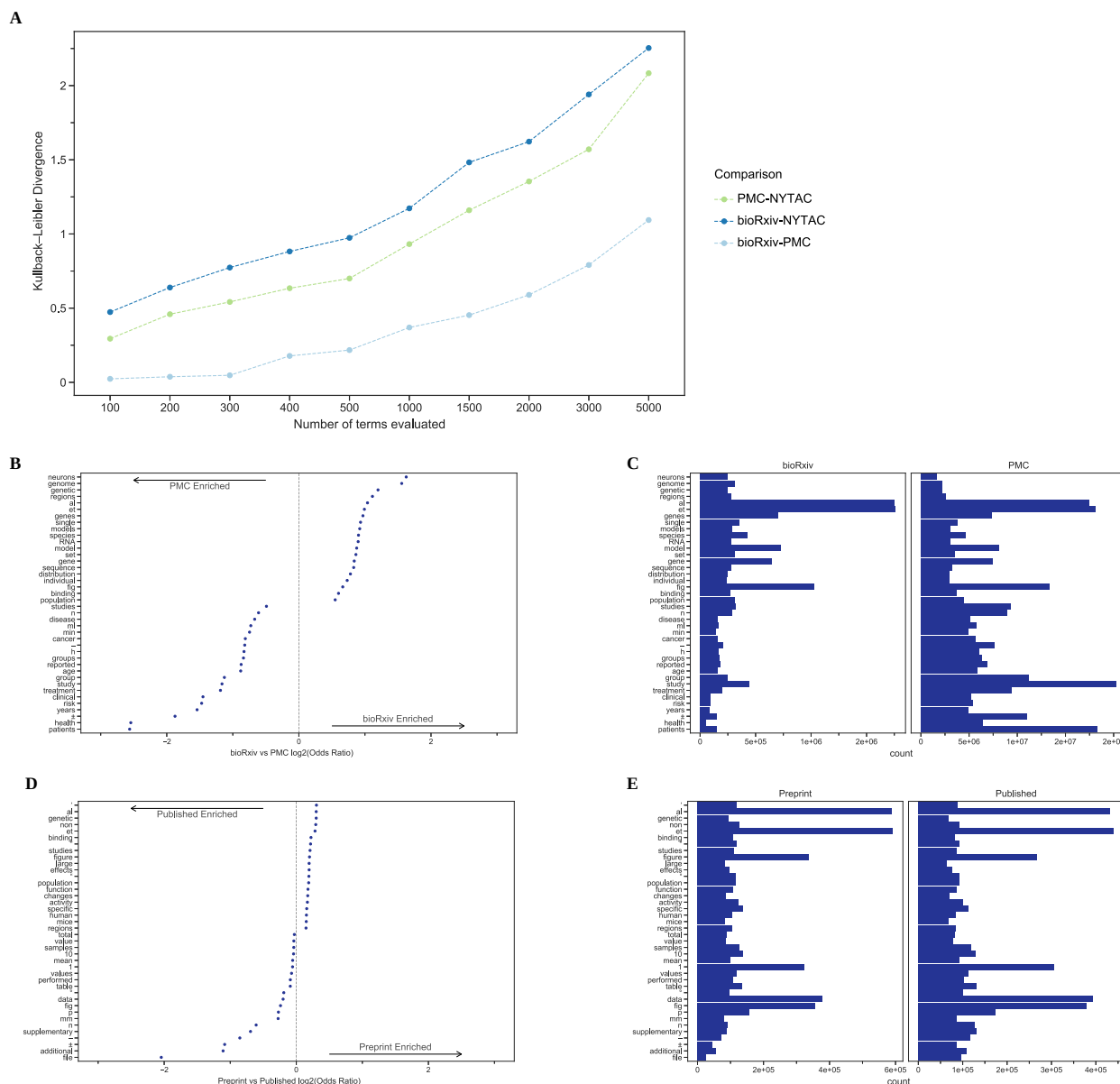
The preprint landscape is rapidly changing, and the number of bioRxiv preprints in our data download (71,118) was nearly double that of a recent study that reported on a snapshot with 37,648 preprints [58]. Because the rate of change is rapid, we first analyzed category data and compared our results with previous findings. As in previous reports [58], neuroscience remains the most common category of preprint followed by bioinformatics (Supplemental Figure S1). Microbiology, which was fifth in the most recent report [58], has now surpassed evolutionary biology and genomics to move into third. When authors upload their preprints, they select from three result category types: new results, confirmatory results or contradictory results. We found that nearly all preprints (97.5%) were categorized as new results, which is consistent with reports on a smaller set [59]. Taken together, the results suggest that while bioRxiv has experienced dramatic growth, the way in which it is being used appears to have remained consistent in recent years.

### Global analysis reveals similarities and differences between bioRxiv and PMC

**Table 1:** Summary statistics for the bioRxiv, PMC, and NYTAC corpora.

Metric	bioRxiv	PMC	NYTAC
document count	71,118	1,977,647	1,855,658
sentence count	22,195,739	480,489,811	72,171,037
token count	420,969,930	8,597,101,167	1,218,673,384
stopword count	158,429,441	3,153,077,263	559,391,073
avg. document length	312.10	242.96	38.89
avg. sentence length	22.71	21.46	19.89
negatives	1,148,382	24,928,801	7,272,401
coordinating conjunctions	14,295,736	307,082,313	38,730,053
coordinating conjunctions%	3.40%	3.57%	3.18%
pronouns	4,604,432	74,994,125	46,712,553
pronouns%	1.09%	0.87%	3.83%
passives	15,012,441	342,407,363	19,472,053
passive%	3.57%	3.98%	1.60%





**Figure 1:** A. The Kullback–Leibler divergence measures the extent to which the distributions, not specific tokens, differ from each other. The token distribution of bioRxiv and PMC corpora is more similar than these biomedical corpora are to the NYTAC one. B. The major differences in token frequencies for the corpora appear to be driven by the fields that have had the highest uptake of bioRxiv, as terms from neuroscience and genomics are relatively more abundant in bioRxiv. Points indicate the  $\log_2(\text{OR})$  for each token and error bars indicate the 95% confidence intervals. C. Of the terms that differ between bioRxiv and PMC, the most abundant in bioRxiv are “et” and “al” while the most abundant in PMC is “study.” D. The major differences in token frequencies for preprints and their corresponding published version often appear to be associated with typesetting and supplementary or additional materials. Points indicate the  $\log_2(\text{OR})$  for each token and error bars indicate the 95% confidence intervals. E. The tokens with the largest absolute differences in abundance appear to be stylistic.

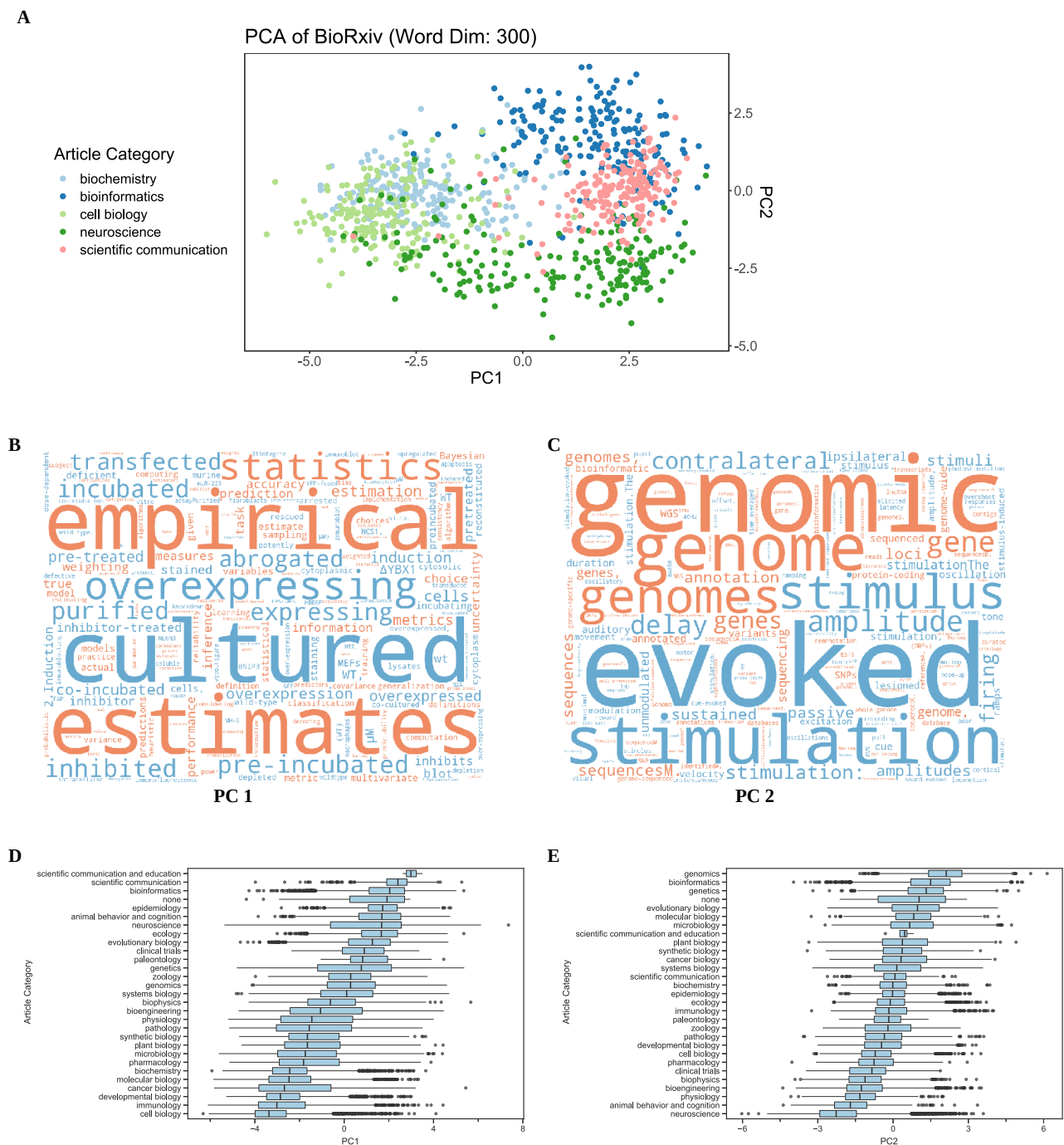
We first compared the overall text of bioRxiv with PMC, adding a corpus of professionally written but non-biomedical text (NYTAC) as a control. Documents on bioRxiv were slightly longer than those on PMC, but both were much longer than those from NYTAC (Table 1). Other than length, both corpora were otherwise quite similar. The average sentence length, fraction of pronouns, and the use of the passive voice were all more similar between bioRxiv and PMC than they were to NYTAC (Table 1). The Kullback–Leibler divergence measures the extent to which two distributions, but not the specific entities that comprise those distributions, differ. The distribution of term frequencies in bioRxiv and PMC was low, especially among the top few hundred tokens (Figure 1A). Differences began to emerge over more terms, but remained much lower than when the biomedical corpora are compared with NYTAC.

Examining the frequencies of individual terms revealed differences between the biomedical corpora. Previous work examining author-selected categories has reported that fields appear to have preprinted unevenly, with certain life sciences research fields having more uptake than others [58]. However, it was possible that authors simply selected certain fields preferentially but that the content was similar to the broader corpus of life sciences text. We directly examined this by comparing term frequencies between bioRxiv and PMC. We found that among the terms that differed the most many were associated with certain life sciences research fields. Terms like “neurons” “genome” and “genetic”, which are common in genomics and neuroscience, were more common in bioRxiv than PMC while others associated with clinical research, such as “clinical” “patients” and “treatment” were more common in PMC (Figure 1B and 1C).

We next controlled for differences in the body of documents to identify term-level changes associated with the publication process itself by examining only pairs of preprints and their corresponding publication (Figure 1D and 1E). The tokens that differed included “et” “al”, “±”, “-” and others that appeared to be typesetting related. Certain changes appeared to be related to journal styles: “figure” was more common in bioRxiv while “fig” was relatively more common in PMC. Other changes appeared to be associated with an increasing reference to content external to the manuscript itself: the tokens “supplementary”, “additional” and “file” were all more common in PMC than bioRxiv suggesting that journals are not simply replacing one token with another but that there are more mentions of such content after peer review.

Taken together these results suggested that the structure of the text in documents on bioRxiv was similar to that on PMC. The differences in uptake across fields are supported not only by differences in authors’ categorization of their articles but also in the text of the articles themselves. At the level of individual manuscripts, the terms that change the most appear to be associated with typesetting, journal style, and an increasing reliance on additional materials after peer review.

## **Document embeddings derived from bioRxiv reveal fields and subfields**



**Figure 2:** A. Principal components (PC) analysis of bioRxiv word2vec embeddings groups documents by author-selected categories. We visualized documents from key categories on a scatterplot for the first two PCs. The first PC separated cell biology from informatics-related fields. The second PC separated bioinformatics from neuroscience. Certain neuroscience papers appeared to be more associated with the cellular biology direction of PC1, while others appeared to be more associated with the informatics-related direction, which suggested that the concepts captured by PCs were not exclusively related to field. B. A word cloud visualization of PC1, which separated informatics disciplines (positive direction) from cell biology (negative direction) showed that tokens “empirical” “estimates” and “statistics” characterized the positive direction while “cultured” and “overexpressing” characterized the negative one. Each word cloud depicts the cosine similarity score between tokens and the second PC. Tokens in orange were most similar to the PC’s positive direction while tokens in blue were most similar to the PC’s negative direction. The size of each token indicates the magnitude of the similarity. C. A word cloud visualization of PC2, which separated bioinformatics from neuroscience, showed that tokens “genomic” “genome” and “genomes” characterized the positive direction while “evoked” “stimulus” and “stimulation” characterized the negative one. D. Examining PC1 values for each article by category created a continuum from informatics-related fields on the top through cell biology on the bottom. Certain article categories (neuroscience, genetics) were spread throughout PC1 values. E. Examining PC2 values for each article by category revealed fields like genomics, bioinformatics, and genetics on the top and neuroscience and behavior on the bottom.

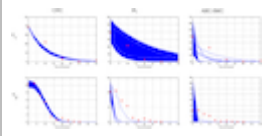
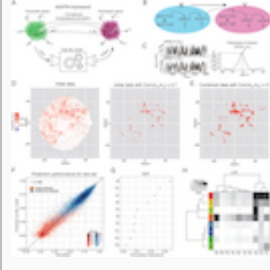

Document embeddings provide a means to categorize the language of documents in a way that takes into account the similarities between terms [33,60,61]. We first trained word embeddings using a 300-dimensional word2vec continuous bag of words model. We combined word embeddings to produce an embedding for each bioRxiv or PMC document by calculating the average of all words present in each respective document. This placed each document in a 300-dimensional space where each individual dimension is arbitrary. To provide more structure to the dataset, we examined the predominant patterns in these embeddings by performing principal components analysis of bioRxiv. The principal components (PCs) are ordered by the proportion of the variance explained. We found that the first two PCs separated articles from different author-selected categories (Figure 2A).

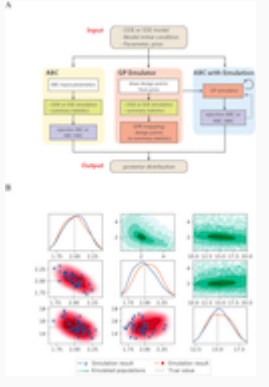
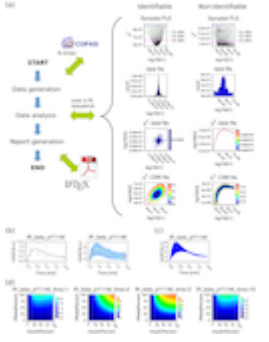

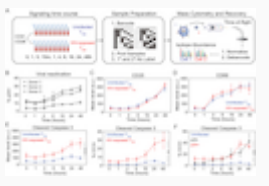
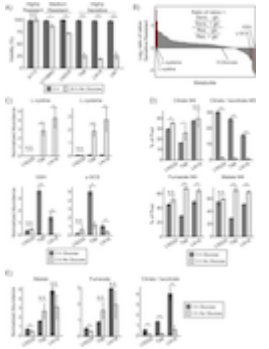
We sought to understand the token patterns that drove these overall differences between documents. We identified the tokens most similar to a PC by calculating the cosine similarity score between tokens' embeddings in the word2vec space and each PC. Visualizing token-PC similarity revealed tokens associated with certain research approaches (Figures 2B and 2C). Examining the value for PC1 across all author-selected categories revealed an ordering of fields from cell biology to informatics-related disciplines (Figure 2D).

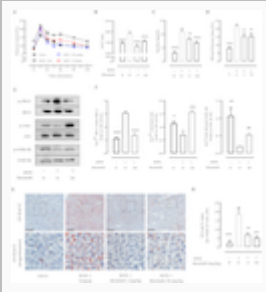
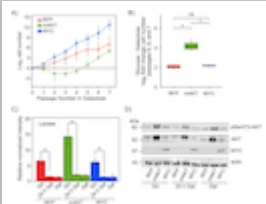
While the PC1 value range for each author-selected category was high, these results suggested that a primary driver in the variability of language use on bioRxiv could be the divide between data science approaches and cell biology ones. A similar analysis for PC2 suggested that neuroscience and genomics present a similar language continuum (Figure 2E). For both of the top two PCs, the submitter-selected category of systems biology preprints was near the middle of the distribution and had a relatively large interquartile range when compared with other categories (Figure 2D and 2E).

We examined the preprints that had the highest and lowest values for PC1 within systems biology (Table {#tbl:five\_pc1\_table}). The preprints with the highest five PC values [62,63,64,65] included software packages, machine learning analyses, and other computational biology manuscripts. The preprints with the lowest five PC values [66,67,68,69,70] were focused on signaling. We provide the top 50 PCs of bioRxiv embeddings within our online repository (see Software and Data Availability).

**Table 2:** PC1 divided the author-selected category of systems biology preprints along an axis from computational to molecular approaches.

Title [citation]	PC1	License	Figure Thumbnail
Conditional Robust Calibration (CRC): a new computational Bayesian methodology for model parameters estimation and identifiability analysis [62]	4.700554908074704	None	
Machine learning of stochastic gene network phenotypes [63]	4.410660604449826	CC-BY-NC-ND	
Notions of similarity for computational biology models [71]	4.355295926618207	CC-BY-NC-ND	

Title [citation]	PC1	License	Figure Thumbnail
GpABC: a Julia package for approximate Bayesian computation with Gaussian process emulation [64]	4.351517618262304	CC-BY-NC-ND	 <p>The figure for GpABC consists of two parts. Part A is a flowchart showing the workflow: 'Input' (Data or GP model, Model input parameters) leads to 'GP Emulation' (GP Emulation, GP Emulation, GP Emulation) and 'ABC with Emulation' (ABC with Emulation, ABC with Emulation). Part B is a grid of plots showing 'Simulation result' (red) and 'True value' (blue) for various parameters, including 'GP Emulation' and 'ABC with Emulation'.</p>
SBpipe: a collection of pipelines for automating repetitive simulation and analysis tasks [65]	4.321847854182741	CC-BY-NC-ND	 <p>The figure for SBpipe shows a flowchart of the pipeline: 'Data generation' (Data generation, Data generation, Data generation) leads to 'Data analysis' (Data analysis, Data analysis, Data analysis) and 'Report generation' (Report generation, Report generation, Report generation). Below the flowchart are several plots showing 'Simulation result' (red) and 'True value' (blue) for various parameters, including 'Data generation', 'Data analysis', and 'Report generation'.</p>
Spatiotemporal proteomics uncovers cathepsin-dependent host cell death during bacterial infection [66]	-4.263964235099807	CC-BY-ND	 <p>The figure for Spatiotemporal proteomics shows a complex diagram of the host cell death pathway, including 'Cathepsin', 'Lysosome', and 'Mitochondria'. Below the diagram are several plots showing 'Simulation result' (red) and 'True value' (blue) for various parameters, including 'Cathepsin', 'Lysosome', and 'Mitochondria'.</p>
Systems analysis by mass cytometry identifies susceptibility of latent HIV-infected T cells to targeting of p38 and mTOR pathways [67]	-4.279016673409032	CC-BY-NC-ND	 <p>The figure for Systems analysis by mass cytometry shows a flowchart of the analysis pipeline: 'Data generation' (Data generation, Data generation, Data generation) leads to 'Data analysis' (Data analysis, Data analysis, Data analysis) and 'Report generation' (Report generation, Report generation, Report generation). Below the flowchart are several plots showing 'Simulation result' (red) and 'True value' (blue) for various parameters, including 'Data generation', 'Data analysis', and 'Report generation'.</p>
NADPH consumption by L-cystine reduction creates a metabolic vulnerability upon glucose deprivation [68]	-4.592209988884236	None	 <p>The figure for NADPH consumption shows a flowchart of the metabolic pathway: 'Data generation' (Data generation, Data generation, Data generation) leads to 'Data analysis' (Data analysis, Data analysis, Data analysis) and 'Report generation' (Report generation, Report generation, Report generation). Below the flowchart are several plots showing 'Simulation result' (red) and 'True value' (blue) for various parameters, including 'Data generation', 'Data analysis', and 'Report generation'.</p>

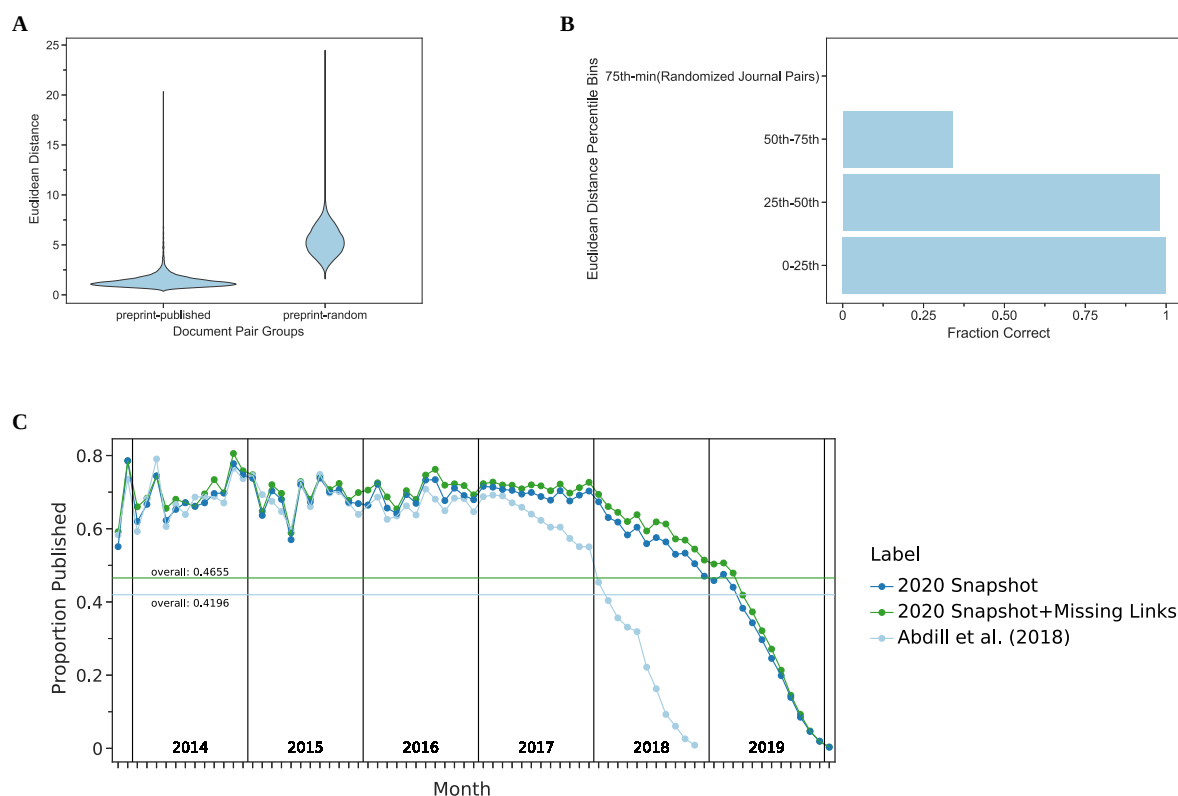
Title [citation]	PC1	License	Figure Thumbnail
Inhibition of Bruton's tyrosine kinase reduces NF-kB and NLRP3 inflammasome activity preventing insulin resistance and microvascular disease [69]	-4.736613689905791	None	
AKT but not MYC promotes reactive oxygen species-mediated cell death in oxidative culture [70]	-4.826793742506695	None	

## Document embedding similarities reveal unannotated preprint-publication pairs

Metaresearch into bioRxiv, including our own, relies on annotations of preprints that have been published to their corresponding peer reviewed publication. Many journals require that authors update preprints with links to the published version of their article. This is accomplished in two ways: bioRxiv may detect the link and automatically add it or authors may notify bioRxiv that their preprint was published. However, bioRxiv establishes these links based on consistency in metadata (i.e., title, author names, etc). Article titles, author lists, and other elements may change as a result of the peer review process. In these cases, if authors do not notify bioRxiv that the preprint has since been published then the scientific record remains incomplete.

Based on our finding that document embeddings captured fields and subfields, we expected that preprint-publication pairs with similar embedding values may represent the same document. We examined the extent to which annotated preprint-publication pairs were closer in this space that preprints were to a random paper published in the same journal in which the preprint was eventually published using already annotated pairs. We found that distances between preprints and their corresponding published versions were nearly always lower than preprints paired with a random article published in the same journal (Figure 3A).





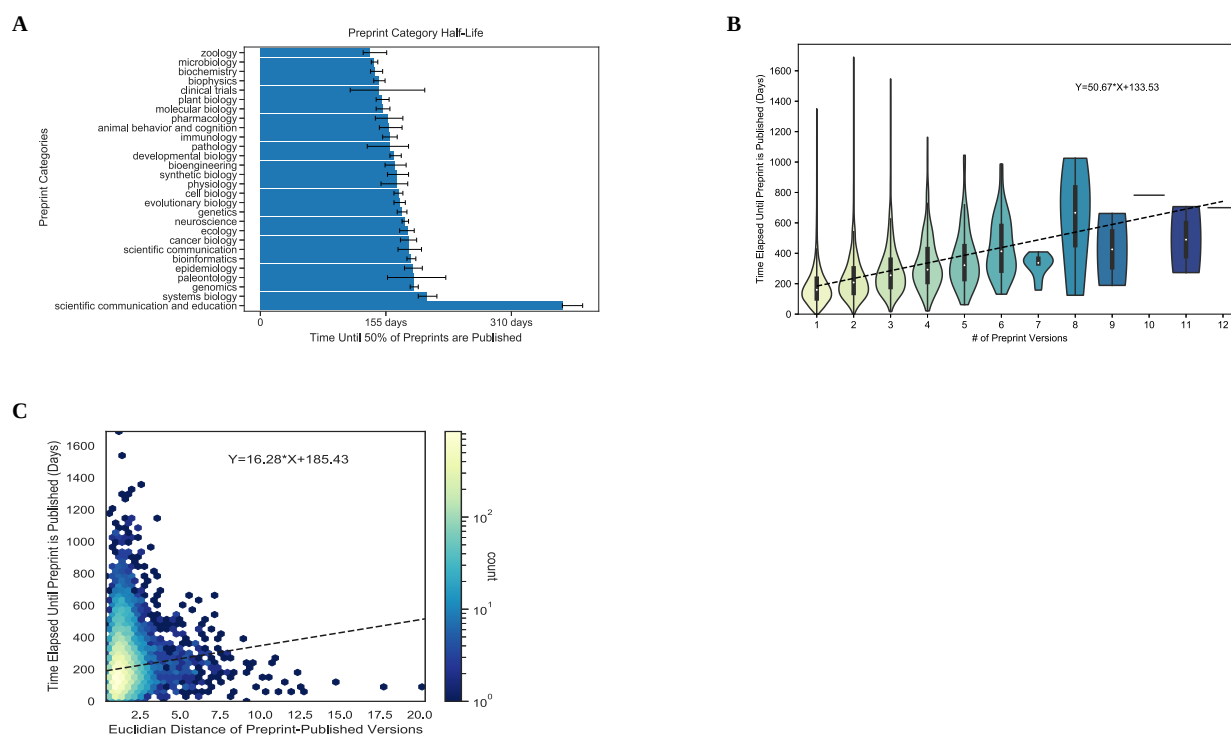
**Figure 3:** A. Preprints are closer in document embedding space to their corresponding peer reviewed publication than they are to random papers published in the same journal. B. Potential preprint-publication pairs that are unannotated but within the 50th percentile of all preprint-publication pairs in the document embedding space are likely represent true preprint-publication pairs. We depict the fraction of true positives over the total number of pairs in each bin. Accuracy is derived from curation of a randomized list of 200 potential pairs (50 per quantile) performed in duplicate with a third rater used in the case of disagreement. C. Most preprints are eventually published. We show the publication rate of preprints since bioRxiv first started. The x-axis represents months since bioRxiv started and the y-axis represents the proportion of preprints published. The light blue line represents the publication rate estimated by Abdill et al. [58]. The dark blue line represents the updated publication rate using only CrossRef-derived annotations while the dark green line includes annotations derived from similarity in the embedding space. The horizontal lines represent the overall proportion of preprints that are were published as of the time of the annotation snapshot.

Based on this finding, we analyzed preprint-publication pairs that were close in document space but not annotated as such. We separated these pairs into four quantiles with the first three based on the distribution of preprint-publication distances and the fourth going from the 75th percentile in the preprint-publication pair space to the smallest value observed for the preprint-random set. We then selected 50 preprint-publication pairs from each of these sets and shuffled them to create a random list of 200 possible pairs. Two scientists then examined the pairs in these randomized lists determining if an preprint-publication pair represented the peer reviewed publication of the primary content described in the preprint. Across the entire list we found a high inter-rater reliability of with a Cohen's Kappa [72] of 0.92. In the case of disagreements, a third scientist more carefully examined the pairs and made a final determination. Of the 200 pairs that we examined, approximately 98% of pairs with an embedding distance in the 0-25th and 25th-50th percentile bins were scored as true matches (Figure 3B). These two bins contained 1,720 preprint-article pairs, suggesting that many preprints have been published but not previously connected with their published versions.

We overlaid these new annotations onto existing annotations to reassess the overall preprint publication rate reported by Abdill et al. [58]. Our filtering criteria were intentionally stringent, so the increased estimate of publication rate amounts to a few percent (Figure 3C). We noticed that there was a particular enrichment of unannotated but published preprints in the 2017-2018 interval. We would expect a higher proportion of such preprints before the year 2019 (many of which may not

have been published yet); however, we did not expect to observe relatively few missed annotations before 2017. It is possible that as the number of preprints grows, it has become harder to establish links. Alternatively, authors now adopting preprinting practices may be less likely to notify the preprint server upon publication if the links are not automatically detected. In any case, future work, especially that which aims to assess the fraction of preprints that are eventually published, should account for the possibility of missed annotations. We supplied our set of 1720 high-confidence annotations to the bioRxiv staff.

## Preprints with more versions or more text changes took longer to publish



**Figure 4:** A. Preprints with more substantial text changes took longer to be published. The x-axis shows the Euclidean distance between document representations of the first version of a preprint and its peer reviewed form. The y-axis shows the number of days elapsed between when the first version of a preprint posted on bioRxiv and the time a preprint is published. The color bar on the right represents the density of each hexbin in this plot: more dense regions are shown in a brighter color. B. Preprints with more versions were associated with a longer time to publish. The x-axis shows the number of versions of a preprint that were posted on bioRxiv. The y-axis shows the number of days that elapsed between when the first version of a preprint was posted on bioRxiv and the date at which the peer reviewed publication appeared. The density of observations are depicted with the violin plot with an embedded boxplot. C. Author-selected categories were associated with modest differences in the time to publish. Categories are shown on the y-axis. The x-axis shows the median time-to-publish for each category. Error bars represent 95% confidence intervals for each preprint category's median time to publication.

The process of peer review includes a number of steps which take variable amounts of time [73]. Comparing bioRxiv preprints with their corresponding published version provides an opportunity to better understand peer review and publishing. We examined how long it took to publish preprints in their peer reviewed form for each author-selected category (Figure 4C). We considered a preprint's endpoint to be publication and used the Kaplan-Meier estimator to estimate how much time has elapsed. Of the most abundant preprint categories microbiology was the fastest to publish (140 days, (137, 145 days) [95% CI]) and genomics was the slowest (190 days, (185, 195 days) [95% CI]). Though we did observe category-specific differences, these differences were generally modest. One exception was the Scientific Communication and Education category, which took substantially longer to be peer reviewed and published (373 days, (373, 398 days) [95% CI]).

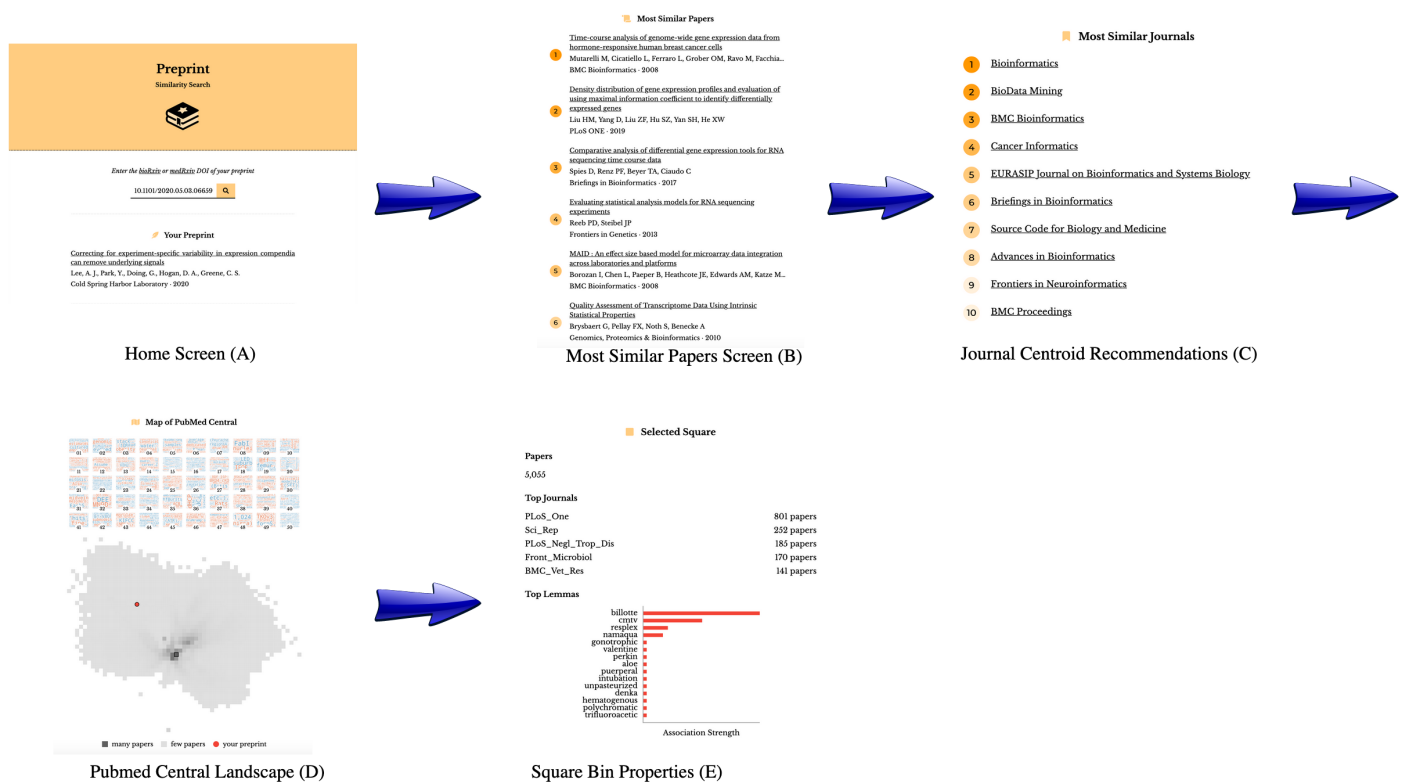
Preprint authors have the opportunity to update their preprint at any point until the peer reviewed publication appears. We examined whether or not the number of versions of a preprint was associated with a change in the time to publish. We found that preprints with more versions generally took longer to publish (Figure 4B). Linear regression comparing the time to publish with the number of versions revealed an increasing relationship, with each version being associated with an additional 51 days before publication. This time period seems broadly compatible with the amount of time it would take to receive reviews and revise a manuscript, suggesting that many authors may be updating their preprints in response to peer reviews or other external feedback.

We next used the full text content to examine the extent to which the magnitude of change to the text between the preprint and peer reviewed publication was associated with a change in the time to publication. We used distance in the document embedding space to quantify the amount of change to the text. We first established a baseline by selecting pairs of random preprints in the author-selected bioinformatics category and found a mean embedding distance between random pairs of 5.068. Preprints that were further in the embedding space from their corresponding peer reviewed publication took longer to publish, with each unit in the embedding space corresponding to approximately sixteen additional days (Figure 4A). The number of version and document distance effects, taken together, support a model where preprints that are reviewed more or that require larger revisions take longer to publish.

## **Preprints with similar document embeddings share publication venues**

We hypothesized that document embedding might capture information associated with the eventual publication venue. We used the annotated preprint-publication pairs as a fully held out test set. We embedded all PMC papers except those annotated to have a corresponding preprint into the bioRxiv embedding space. For journals, we averaged the document representation of all papers in the journal to produce a journal centroid. We used a k-nearest neighbor classifier trained on the remaining PMC documents to identify the most similar published papers and assessed performance - the extent to which the journal a paper appeared in the returned list of ten neighbors - through cross validation and the held out test set (Supplemental Figure S4). We found an enrichment of papers published within the same journal among the neighbors of the target paper.

We developed an online app that returns the nearest published papers and journals that are nearest to a preprint. Users supply digital object identifiers (DOIs) from bioRxiv or medRxiv. The application downloads the article from the preprint repository, converts the PDF to text, calculates a document embedding score, and returns the ten papers and journals with the most similar representations in the embedding space (Figure 5). We also sought to display the position of each preprint in the overall landscape. We required an embedding that could be rapidly calculated for new documents, so we used the Sparse Autoencoder for Unsupervised Clustering, Imputation, and Embedding (SAUCIE) approach previously described for the analysis of single-cell gene expression data [74]. The user-requested preprint's location in this space is then displayed and users can select regions to identify the terms most associated with those regions. Users can also explore the terms associated with the top 50 PCs derived from the document embeddings and those PCs vary across the document landscape.



**Figure 5:** The preprint similarity search app workflow allows users to examine where an individual preprint falls in the overall document landscape. A. Starting with the home screen, users can paste in a bioRxiv or medRxiv DOI, which sends a request to bioRxiv or medRxiv. Next the app preprocesses the requested preprint and returns a listing of (B) the top ten most similar papers and (C) the ten closest journals. D. The app also displays the location of the query preprint in PMC. E. Users can select a square within the landscape to examine statistics associated with the square including the top journals by article count in that square and the odds ratio of tokens.

## Conclusions

We analyzed the full text of bioRxiv, comparing it with the open access portion of PMC, and found that the overall manner of writing is consistent with the biomedical literature. A benefit of analyzing bioRxiv text is that we can compare preprint and published versions to examine the influence of the publication process. Token-level analyses suggest that differences between corpora are driven by fields, while comparisons of preprints with their corresponding publication reveals differences in typesetting and supplementary materials.

Previous analyses have focused on article metadata, including which papers are heavily downloaded and discussed [58] and by which communities discuss them [75]. We found that some preprints are highly similar to published articles within the PMC open access corpus, and a detailed examination revealed that many preprints were published but not previously annotated. Using the full text of documents to correct these missing annotations provides a comprehensive understanding of the extent to which preprints are published and this correction resulted in a publication rate that is higher than previously estimated. Importantly, this only accounts for papers that are published open access, so our analysis should be considered to raise the lower bound but the truly published fraction is likely to be higher.

Our work presents a first step towards understanding the process by which peer review alters life sciences papers through an analysis of the full text of preprints and their corresponding publication. Our broad-based examination suggests certain changes, but the scale of changes appears modest. We lay the groundwork for future work that aims to identify sentences and claims that are altered during review. Our finding that document embeddings reveal manuscripts with similar outcomes is likely to be important to new tools that accelerate publishing, including those that automatically perform integrity checks and other critically important contributions. We also found that document

embeddings captured many elements of how authors use language in preprints. Fields were separated first on an informatics / molecular axis, and author-selected categories could be distinguished based on embeddings. Document embeddings were also associated with the eventual journal at which the work was published. Based on this observation, we supply a web application that returns the papers and journals that are most similar to a bioRxiv or medRxiv preprint.

## Software and Data Availability

---

An online version of this manuscript is available under a Creative Commons Attribution License at [https://greenelab.github.io/annorxiver\\_manuscript/](https://greenelab.github.io/annorxiver_manuscript/). Source for the research portions of this project is dual licensed under the BSD 3-Clause and Creative Commons Public Domain Dedication Licenses at <https://github.com/greenelab/annorxiver>. The journal recommendation website can be found at <https://greenelab.github.io/annorxiver-journal-recommender/>, and code for the website is available under a BSD-2-Clause Plus Patent License at <https://github.com/greenelab/preprint-similarity-search>. Full text access for the bioRxiv repository is available at <https://www.biorxiv.org/tdm>. Access to PubMed Central's Open Access subset is available on NCBI's FTP server at <https://www.ncbi.nlm.nih.gov/pmc/tools/ftp/>. Access to the New York Times Annotated Corpus (NYTAC) is available upon request with the Linguistic Data Consortium at <https://catalog.ldc.upenn.edu/LDC2008T19>.

## Acknowledgments

---

The authors would like to thank Ariel Hippen Anderson for evaluating potential missing preprint to published version links. We also would like to thank Richard Sever and the bioRxiv team for their assistance with access to and support with questions about preprint full text downloaded from bioRxiv.

## Funding

---

This work was supported by grants from the Gordon Betty Moore Foundation (GBMF4552) and the National Institutes of Health's National Human Genome Research Institute (NHGRI) under awards T32 HG00046 and R01 HG010067.

## Competing Interests

---

Marvin Thielk receives a salary from Elsevier Inc. where he contributes NLP expertise to health content operations. Elsevier and Marvin provided no restrictions on results or interpretations that could be published in this manuscript. The opinions expressed here do not reflect the official policy or positions held by Elsevier Inc.

# References

---

**1. Scientific communication pathways: an overview and introduction to a symposium**

David F. Zaye, W. V. Metanovski

*Journal of Chemical Information and Computer Sciences* (2002-05-01) <https://doi.org/bwsxhg>

DOI: [10.1021/ci00050a001](https://doi.org/10.1021/ci00050a001)

**2. The Transition from Paper to Electronic Journals**

Hak Joon Kim

*The Serials Librarian* (2001-11-19) <https://doi.org/d7rnh2>

DOI: [10.1300/j123v41n01\\_04](https://doi.org/10.1300/j123v41n01_04)

**3. Preprints: What Role Do These Have in Communicating Scientific Results?**

Susan A. Elmore

*Toxicologic Pathology* (2018-04-08) <https://doi.org/ghdd7c>

DOI: [10.1177/0192623318767322](https://doi.org/10.1177/0192623318767322) · PMID: [29628000](https://pubmed.ncbi.nlm.nih.gov/29628000/) · PMCID: [PMC5999550](https://pubmed.ncbi.nlm.nih.gov/PMC5999550/)

**4. The prehistory of biology preprints: A forgotten experiment from the 1960s**

Matthew Cobb

*PLOS Biology* (2017-11-16) <https://doi.org/c6ww>

DOI: [10.1371/journal.pbio.2003995](https://doi.org/10.1371/journal.pbio.2003995) · PMID: [29145518](https://pubmed.ncbi.nlm.nih.gov/29145518/) · PMCID: [PMC5690419](https://pubmed.ncbi.nlm.nih.gov/PMC5690419/)

**5. arXiv.org: the Los Alamos National Laboratory e-print server**

Gerry McKiernan

*International Journal on Grey Literature* (2000-09) <https://doi.org/fg8pw7>

DOI: [10.1108/14666180010345564](https://doi.org/10.1108/14666180010345564)

**6. bioRxiv: the preprint server for biology**

Richard Sever, Ted Roeder, Samantha Hindle, Linda Sussman, Kevin-John Black, Janet Argentine, Wayne Manos, John R. Inglis

*Cold Spring Harbor Laboratory* (2019-11-06) <https://doi.org/ggc46z>

DOI: [10.1101/833400](https://doi.org/10.1101/833400)

**7. medRxiv.org - the preprint server for Health Sciences** <https://www.medrxiv.org/>

**8. The Second Wave of Preprint Servers: How Can Publishers Keep Afloat?**

By

*The Scholarly Kitchen* (2019-10-16) <https://scholarlykitchen.sspnet.org/2019/10/16/the-second-wave-of-preprint-servers-how-can-publishers-keep-afloat/>

**9. Rxivist.org: Sorting biology preprints using social media and readership metrics**

Richard J. Abdill, Ran Blekhman

*PLOS Biology* (2019-05-21) <https://doi.org/dm27>

DOI: [10.1371/journal.pbio.3000269](https://doi.org/10.1371/journal.pbio.3000269) · PMID: [31112533](https://pubmed.ncbi.nlm.nih.gov/31112533/) · PMCID: [PMC6546241](https://pubmed.ncbi.nlm.nih.gov/PMC6546241/)

**10. How the Scientific Community Reacts to Newly Submitted Preprints: Article Downloads, Twitter Mentions, and Citations**

Xin Shuai, Alberto Pepe, Johan Bollen

*PLoS ONE* (2012-11-01) <https://doi.org/f4cw6r>

DOI: [10.1371/journal.pone.0047523](https://doi.org/10.1371/journal.pone.0047523) · PMID: [23133597](https://pubmed.ncbi.nlm.nih.gov/23133597/) · PMCID: [PMC3486871](https://pubmed.ncbi.nlm.nih.gov/PMC3486871/)



11. **Quantifying and contextualizing the impact of bioRxiv preprints through automated social media audience segmentation**  
Jedidiah Carlson, Kelley Harris  
*Cold Spring Harbor Laboratory* (2020-03-10) <https://doi.org/ghdd66>  
DOI: [10.1101/2020.03.06.981589](https://doi.org/10.1101/2020.03.06.981589)
12. **Abstract**  
eLife Sciences Publications, Ltd  
(2019-05-09) <https://doi.org/gf5cqt>  
DOI: [10.7554/elife.45133.001](https://doi.org/10.7554/elife.45133.001)
13. **An analysis of published journals for papers posted on bioRxiv**  
Hiroyuki Tsunoda, Yuan Sun, Masaki Nishizawa, Xiaomin Liu, Kou Amano  
*Proceedings of the Association for Information Science and Technology* (2019-10-18)  
<https://doi.org/ggz7f9>  
DOI: [10.1002/pr2.175](https://doi.org/10.1002/pr2.175)
14. **The relationship between bioRxiv preprints, citations and altmetrics**  
Nicholas Fraser, Fakhri Momeni, Philipp Mayr, Isabella Peters  
*Quantitative Science Studies* (2020-04-01) <https://doi.org/gg2cz3>  
DOI: [10.1162/qss\\_a\\_00043](https://doi.org/10.1162/qss_a_00043)
15. **Releasing a preprint is associated with more attention and citations for the peer-reviewed article**  
Darwin Y Fu, Jacob J Hughey  
*eLife* (2019-12-06) <https://doi.org/ghd3mv>  
DOI: [10.7554/elife.52646](https://doi.org/10.7554/elife.52646) · PMID: [31808742](https://pubmed.ncbi.nlm.nih.gov/31808742/) · PMCID: [PMC6914335](https://pubmed.ncbi.nlm.nih.gov/PMC6914335/)
16. **Preprints and Scholarly Communication: An Exploratory Qualitative Study of Adoption, Practices, Drivers and Barriers**  
Andrea Chiarelli, Rob Johnson, Stephen Pinfield, Emma Richens  
*F1000Research* (2019-11-25) <https://doi.org/ghp38z>  
DOI: [10.12688/f1000research.19619.2](https://doi.org/10.12688/f1000research.19619.2) · PMID: [32055396](https://pubmed.ncbi.nlm.nih.gov/32055396/) · PMCID: [PMC6961415](https://pubmed.ncbi.nlm.nih.gov/PMC6961415/)
17. **The Need for Speed: How Quickly Do Preprints Become Published Articles?**  
Rachel Herbert, Kate Gasson, Alex Ponsford  
*SSRN Electronic Journal* (2019) <https://doi.org/ghd3mt>  
DOI: [10.2139/ssrn.3455146](https://doi.org/10.2139/ssrn.3455146)
18. **Technical and social issues influencing the adoption of preprints in the life sciences**  
Naomi C. Penfold, Jessica K. Polka  
*PLOS Genetics* (2020-04-20) <https://doi.org/dtt2>  
DOI: [10.1371/journal.pgen.1008565](https://doi.org/10.1371/journal.pgen.1008565) · PMID: [32310942](https://pubmed.ncbi.nlm.nih.gov/32310942/) · PMCID: [PMC7170218](https://pubmed.ncbi.nlm.nih.gov/PMC7170218/)
19. **Biologists urged to hug a preprint**  
Ewen Callaway, Kendall Powell  
*Nature* (2016-02-16) <https://doi.org/ghdd62>  
DOI: [10.1038/530265a](https://doi.org/10.1038/530265a) · PMID: [26887471](https://pubmed.ncbi.nlm.nih.gov/26887471/)
20. **On the value of preprints: An early career researcher perspective**  
Sarvenaz Sarabipour, Humberto J. Debat, Edward Emmott, Steven J. Burgess, Benjamin Schwessinger, Zach Hensel

*PLOS Biology* (2019-02-21) <https://doi.org/gfw8hd>  
DOI: [10.1371/journal.pbio.3000151](https://doi.org/10.1371/journal.pbio.3000151) · PMID: [30789895](https://pubmed.ncbi.nlm.nih.gov/30789895/) · PMCID: [PMC6400415](https://pubmed.ncbi.nlm.nih.gov/PMC6400415/)

**21. Prepublication Communication of Research Results**

Michael J. Adams, Reid N. Harris, Evan H. C. Grant, Matthew J. Gray, M. Camille Hopkins, Samuel A. Iverson, Robert Likens, Mark Mandica, Deanna H. Olson, Alex Shepack, Hardin Waddle  
*EcoHealth* (2018-08-07) <https://doi.org/ghn66s>  
DOI: [10.1007/s10393-018-1352-3](https://doi.org/10.1007/s10393-018-1352-3) · PMID: [30088185](https://pubmed.ncbi.nlm.nih.gov/30088185/) · PMCID: [PMC6245104](https://pubmed.ncbi.nlm.nih.gov/PMC6245104/)

**22. Peer Review and bioRxiv**

Leslie M. Loew  
*Biophysical Journal* (2016-08) <https://doi.org/ghdd6x>  
DOI: [10.1016/j.bpj.2016.06.035](https://doi.org/10.1016/j.bpj.2016.06.035) · PMID: [27508451](https://pubmed.ncbi.nlm.nih.gov/27508451/) · PMCID: [PMC4982934](https://pubmed.ncbi.nlm.nih.gov/PMC4982934/)

**23. Textual Analysis in Accounting and Finance: A Survey**

TIM LOUGHRAN, BILL MCDONALD  
*Journal of Accounting Research* (2016-09) <https://doi.org/gc3hf7>  
DOI: [10.1111/1475-679x.12123](https://doi.org/10.1111/1475-679x.12123)

**24. The textual characteristics of traditional and Open Access scientific journals are similar**

Karin Verspoor, K Bretonnel Cohen, Lawrence Hunter  
*BMC Bioinformatics* (2009-06-15) <https://doi.org/b973tn>  
DOI: [10.1186/1471-2105-10-183](https://doi.org/10.1186/1471-2105-10-183) · PMID: [19527520](https://pubmed.ncbi.nlm.nih.gov/19527520/) · PMCID: [PMC2714574](https://pubmed.ncbi.nlm.nih.gov/PMC2714574/)

**25. Current findings from research on structured abstracts**

James Hartley  
*Journal of the Medical Library Association : JMLA* (2004-07)  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC442180/>  
PMID: [15243644](https://pubmed.ncbi.nlm.nih.gov/15243644/) · PMCID: [PMC442180](https://pubmed.ncbi.nlm.nih.gov/PMC442180/)

**26. A survey on annotation tools for the biomedical literature**

M. Neves, U. Leser  
*Briefings in Bioinformatics* (2012-12-18) <https://doi.org/f5vzsj>  
DOI: [10.1093/bib/bbs084](https://doi.org/10.1093/bib/bbs084) · PMID: [23255168](https://pubmed.ncbi.nlm.nih.gov/23255168/)

**27. PubTator central: automated concept annotation for biomedical full text articles**

Chih-Hsuan Wei, Alexis Allot, Robert Leaman, Zhiyong Lu  
*Nucleic Acids Research* (2019-07-02) <https://doi.org/ggzfsc>  
DOI: [10.1093/nar/gkz389](https://doi.org/10.1093/nar/gkz389) · PMID: [31114887](https://pubmed.ncbi.nlm.nih.gov/31114887/) · PMCID: [PMC6602571](https://pubmed.ncbi.nlm.nih.gov/PMC6602571/)

**28. Coreference annotation and resolution in the Colorado Richly Annotated Full Text (CRAFT) corpus of biomedical journal articles**

K. Bretonnel Cohen, Arrick Lanfranchi, Miji Joo-young Choi, Michael Bada, William A. Baumgartner, Natalya Panteleyeva, Karin Verspoor, Martha Palmer, Lawrence E. Hunter  
*BMC Bioinformatics* (2017-08-17) <https://doi.org/ghmbw2>  
DOI: [10.1186/s12859-017-1775-9](https://doi.org/10.1186/s12859-017-1775-9) · PMID: [28818042](https://pubmed.ncbi.nlm.nih.gov/28818042/) · PMCID: [PMC5561560](https://pubmed.ncbi.nlm.nih.gov/PMC5561560/)

**29. The structural and content aspects of abstracts versus bodies of full text journal articles are different**

K Bretonnel Cohen, Helen L Johnson, Karin Verspoor, Christophe Roeder, Lawrence E Hunter  
*BMC Bioinformatics* (2010-09-29) <https://doi.org/b9f6rn>  
DOI: [10.1186/1471-2105-11-492](https://doi.org/10.1186/1471-2105-11-492) · PMID: [20920264](https://pubmed.ncbi.nlm.nih.gov/20920264/) · PMCID: [PMC3098079](https://pubmed.ncbi.nlm.nih.gov/PMC3098079/)

30. **A corpus of full-text journal articles is a robust evaluation tool for revealing differences in performance of biomedical natural language processing tools**  
Karin Verspoor, Kevin Bretonnel Cohen, Arrick Lanfranchi, Colin Warner, Helen L Johnson, Christophe Roeder, Jinho D Choi, Christopher Funk, Yuriy Malenkiy, Miriam Eckert, ... Lawrence E Hunter  
*BMC Bioinformatics* (2012-08-17) <https://doi.org/gb8t7v>  
DOI: [10.1186/1471-2105-13-207](https://doi.org/10.1186/1471-2105-13-207) · PMID: [22901054](https://pubmed.ncbi.nlm.nih.gov/22901054/) · PMCID: [PMC3483229](https://pubmed.ncbi.nlm.nih.gov/PMC3483229/)
31. **From POS tagging to dependency parsing for biomedical event extraction**  
Dat Quoc Nguyen, Karin Verspoor  
*BMC Bioinformatics* (2019-02-12) <https://doi.org/ggsrkw>  
DOI: [10.1186/s12859-019-2604-0](https://doi.org/10.1186/s12859-019-2604-0) · PMID: [30755172](https://pubmed.ncbi.nlm.nih.gov/30755172/) · PMCID: [PMC6373122](https://pubmed.ncbi.nlm.nih.gov/PMC6373122/)
32. **Efficient Estimation of Word Representations in Vector Space**  
Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean  
*arXiv* (2013-09-10) <https://arxiv.org/abs/1301.3781>
33. **Distributed Representations of Sentences and Documents**  
Quoc V. Le, Tomas Mikolov  
*arXiv* (2014-05-26) <https://arxiv.org/abs/1405.4053>
34. **Machine access and text/data mining resources | bioRxiv** <https://www.biorxiv.org/tdm>
35. **PubMed Central: The GenBank of the published literature**  
R. J. Roberts  
*Proceedings of the National Academy of Sciences* (2001-01-16) <https://doi.org/bbn9k8>  
DOI: [10.1073/pnas.98.2.381](https://doi.org/10.1073/pnas.98.2.381) · PMID: [11209037](https://pubmed.ncbi.nlm.nih.gov/11209037/) · PMCID: [PMC33354](https://pubmed.ncbi.nlm.nih.gov/PMC33354/)
36. **Gold open access: the best of both worlds**  
M. A. G. van der Heyden, T. A. B. van Veen  
*Netherlands Heart Journal* (2017-12-01) <https://doi.org/ggzfr9>  
DOI: [10.1007/s12471-017-1064-2](https://doi.org/10.1007/s12471-017-1064-2) · PMID: [29196877](https://pubmed.ncbi.nlm.nih.gov/29196877/) · PMCID: [PMC5758455](https://pubmed.ncbi.nlm.nih.gov/PMC5758455/)
37. **How Papers Get Into PMC** <https://www.ncbi.nlm.nih.gov/pmc/about/submission-methods/>
38. **8.2.2 NIH Public Access Policy**  
[https://grants.nih.gov/grants/policy/nihgps/html5/section\\_8/8.2.2\\_nih\\_public\\_access\\_policy.htm](https://grants.nih.gov/grants/policy/nihgps/html5/section_8/8.2.2_nih_public_access_policy.htm)
39. **PMC Overview** <https://www.ncbi.nlm.nih.gov/pmc/about/intro/>
40. **PMC text mining subset in BioC: about three million full-text articles and growing**  
Donald C Comeau, Chih-Hsuan Wei, Rezarta Islamaj Doğan, Zhiyong Lu  
*Bioinformatics* (2019-09-15) <https://doi.org/ggzfsb>  
DOI: [10.1093/bioinformatics/btz070](https://doi.org/10.1093/bioinformatics/btz070) · PMID: [30715220](https://pubmed.ncbi.nlm.nih.gov/30715220/) · PMCID: [PMC6748740](https://pubmed.ncbi.nlm.nih.gov/PMC6748740/)
41. **Author Manuscripts in PMC** <https://www.ncbi.nlm.nih.gov/pmc/about/authorms/>
42. **The new york times annotated corpus**  
Evan Sandhaus  
*Linguistic Data Consortium, Philadelphia* (2008)
43. **CrossRef Text and Data Mining Services**  
Rachael Lammey

*Insights the UKSG journal* (2015-07-07) <https://doi.org/gg4hp9>  
DOI: [10.1629/uksg.233](https://doi.org/10.1629/uksg.233)

44. **spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing**  
Matthew Honnibal, Ines Montani  
(2017)
45. **Odds Ratio**  
Steven Tenny, Mary R. Hoffman  
*StatPearls* (2020) <http://www.ncbi.nlm.nih.gov/books/NBK431098/>
46. **Software Framework for Topic Modelling with Large Corpora**  
Radim Řehůřek, Petr Sojka  
*Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks* (2010-05-22)
47. **Probabilistic Principal Component Analysis**  
Michael E. Tipping, Christopher M. Bishop  
*Journal of the Royal Statistical Society: Series B (Statistical Methodology)* (1999-08)  
<https://doi.org/b3hjw7>  
DOI: [10.1111/1467-9868.00196](https://doi.org/10.1111/1467-9868.00196)
48. **Scikit-learn: Machine learning in Python**  
F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, ... E. Duchesnay  
*Journal of Machine Learning Research* (2011)
49. **Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions**  
Nathan Halko, Per-Gunnar Martinsson, Joel A. Tropp  
*arXiv* (2014-04-29) <https://arxiv.org/abs/0909.4061>
50. **The *Drosophila* Cortactin Binding Protein 2 homolog, Nausicaa, regulates lamellipodial actin dynamics in a Cortactin-dependent manner**  
Meghan E. O'Connell, Divya Sridharan, Tristan Driscoll, Ipsita Krishnamurthy, Wick G. Perry, Derek A. Applewhite  
*Cold Spring Harbor Laboratory* (2018-07-24) <https://doi.org/gg4hp7>  
DOI: [10.1101/376665](https://doi.org/10.1101/376665)
51. **The *Drosophila* protein, Nausicaa, regulates lamellipodial actin dynamics in a Cortactin-dependent manner**  
Meghan E. O'Connell, Divya Sridharan, Tristan Driscoll, Ipsita Krishnamurthy, Wick G. Perry, Derek A. Applewhite  
*Biology Open* (2019-06-15) <https://doi.org/gg4hp8>  
DOI: [10.1242/bio.038232](https://doi.org/10.1242/bio.038232) · PMID: [31164339](https://pubmed.ncbi.nlm.nih.gov/31164339/) · PMCID: [PMC6602326](https://pubmed.ncbi.nlm.nih.gov/PMC6602326/)
52. **Medium – Where good ideas find you.**  
Medium  
<https://medium.com>
53. **Understanding survival analysis: Kaplan-Meier estimate**  
Jugal Kishore, ManishKumar Goel, Pardeep Khanna

*International Journal of Ayurveda Research* (2010) <https://doi.org/fdft75>  
DOI: [10.4103/0974-7788.76794](https://doi.org/10.4103/0974-7788.76794) · PMID: [21455458](https://pubmed.ncbi.nlm.nih.gov/21455458/) · PMCID: [PMC3059453](https://pubmed.ncbi.nlm.nih.gov/PMC3059453/)

**54. CamDavidsonPilon/lifelines: v0.25.6**

Cameron Davidson-Pilon, Jonas Kalderstam, Noah Jacobson, Sean-Reed, Ben Kuhn, Paul Zivich, Mike Williamson, Abdealijk, Deepyaman Datta, Andrew Fiore-Gartland, ... Jlim13  
*Zenodo* (2020-10-26) <https://doi.org/ghh2d3>  
DOI: [10.5281/zenodo.4136578](https://doi.org/10.5281/zenodo.4136578)

**55. Scikit-learn: Machine Learning in Python**

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Andreas Müller, Joel Nothman, Gilles Louppe, ... Édouard Duchesnay  
*arXiv* (2018-06-06) <https://arxiv.org/abs/1201.0490>

**56. Welcome to pdfminer.six's documentation! — pdfminer.six 20201018 documentation**  
<https://pdfminersix.readthedocs.io/en/latest/index.html>

**57. Assessing the Heterogeneity of Cardiac Non-myocytes and the Effect of Cell Culture with Integrative Single Cell Analysis**

Brian S. Iskra, Logan Davis, Henry E. Miller, Yu-Chiao Chiu, Alexander R. Bishop, Yidong Chen, Gregory J. Aune  
*Cold Spring Harbor Laboratory* (2020-03-05) <https://doi.org/gg9353>  
DOI: [10.1101/2020.03.04.975177](https://doi.org/10.1101/2020.03.04.975177)

**58. Tracking the popularity and outcomes of all bioRxiv preprints**

Richard J Abdill, Ran Blekhan  
*eLife* (2019-04-24) <https://doi.org/gf2str>  
DOI: [10.7554/elife.45133](https://doi.org/10.7554/elife.45133) · PMID: [31017570](https://pubmed.ncbi.nlm.nih.gov/31017570/) · PMCID: [PMC6510536](https://pubmed.ncbi.nlm.nih.gov/PMC6510536/)

**59. Altmetric Scores, Citations, and Publication of Studies Posted as Preprints**

Stylianios Serghiou, John P. A. Ioannidis  
*JAMA* (2018-01-23) <https://doi.org/gftc69>  
DOI: [10.1001/jama.2017.21168](https://doi.org/10.1001/jama.2017.21168) · PMID: [29362788](https://pubmed.ncbi.nlm.nih.gov/29362788/) · PMCID: [PMC5833561](https://pubmed.ncbi.nlm.nih.gov/PMC5833561/)

**60. Efficient Vector Representation for Documents through Corruption**

Minmin Chen  
*arXiv* (2017-07-11) <https://arxiv.org/abs/1707.02377>

**61. Document Network Projection in Pretrained Word Embedding Space**

Antoine Gourru, Adrien Guille, Julien Velcin, Julien Jacques  
*arXiv* (2020-01-17) <https://arxiv.org/abs/2001.05727>

**62. Conditional Robust Calibration (CRC): a new computational Bayesian methodology for model parameters estimation and identifiability analysis**

Fortunato Bianconi, Chiara Antonini, Lorenzo Tomassoni, Paolo Valigi  
*Cold Spring Harbor Laboratory* (2017-10-02) <https://doi.org/gg9393>  
DOI: [10.1101/197400](https://doi.org/10.1101/197400)

**63. Machine learning of stochastic gene network phenotypes**

Kyemyung Park, Thorsten Prüstel, Yong Lu, John S. Tsang  
*Cold Spring Harbor Laboratory* (2019-10-31) <https://doi.org/gg94bm>  
DOI: [10.1101/825943](https://doi.org/10.1101/825943)

64. **GpABC: a Julia package for approximate Bayesian computation with Gaussian process emulation**  
Evgeny Tankhilevich, Jonathan Ish-Horowicz, Tara Hameed, Elisabeth Roesch, Istvan Kleijn, Michael PH Stumpf, Fei He  
*Cold Spring Harbor Laboratory* (2019-09-18) <https://doi.org/gg94bj>  
DOI: [10.1101/769299](https://doi.org/10.1101/769299)
65. **SBpipe: a collection of pipelines for automating repetitive simulation and analysis tasks**  
Piero Dalle Pezze, Nicolas Le Novère  
*Cold Spring Harbor Laboratory* (2017-02-09) <https://doi.org/gg9392>  
DOI: [10.1101/107250](https://doi.org/10.1101/107250)
66. **Spatiotemporal proteomics uncovers cathepsin-dependent host cell death during bacterial infection**  
Joel Selkig, Nan Li, Jacob Bobonis, Annika Hausmann, Anna Sueki, Haruna Imamura, Bachir El Debs, Gianluca Sigismondo, Bogdan I. Florea, Herman S. Overkleeft, ... Athanasios Typas  
*Cold Spring Harbor Laboratory* (2018-11-07) <https://doi.org/gg94bc>  
DOI: [10.1101/455048](https://doi.org/10.1101/455048)
67. **Systems analysis by mass cytometry identifies susceptibility of latent HIV-infected T cells to targeting of p38 and mTOR pathways**  
Linda E. Fong, Victor L. Bass, Serena Spudich, Kathryn Miller-Jensen  
*Cold Spring Harbor Laboratory* (2018-07-19) <https://doi.org/gg9398>  
DOI: [10.1101/371922](https://doi.org/10.1101/371922)
68. **NADPH consumption by L-cystine reduction creates a metabolic vulnerability upon glucose deprivation**  
James H. Joly, Alireza Delfarah, Philip S. Phung, Sydney Parrish, Nicholas A. Graham  
*Cold Spring Harbor Laboratory* (2019-08-13) <https://doi.org/gg94bf>  
DOI: [10.1101/733162](https://doi.org/10.1101/733162)
69. **Inhibition of Bruton's tyrosine kinase reduces NF- $\kappa$ B and NLRP3 inflammasome activity preventing insulin resistance and microvascular disease**  
Gareth S. D. Purvis, Massimo Collino, Haidee M. A. Tavio, Fausto Chiazza, Caroline E. O'Riodan, Lynda Zeboudj, Nick Guisot, Peter Bunyard, David R. Greaves, Christoph Thiemermann  
*Cold Spring Harbor Laboratory* (2019-08-28) <https://doi.org/gg94bg>  
DOI: [10.1101/745943](https://doi.org/10.1101/745943)
70. **AKT but not MYC promotes reactive oxygen species-mediated cell death in oxidative culture**  
Dongqing Zheng, Jonathan H. Sussman, Matthew P. Jeon, Sydney T. Parrish, Alireza Delfarah, Nicholas A. Graham  
*Cold Spring Harbor Laboratory* (2019-09-01) <https://doi.org/gg94bh>  
DOI: [10.1101/754572](https://doi.org/10.1101/754572)
71. **Notions of similarity for computational biology models**  
Ron Henkel, Robert Hoehndorf, Tim Kacprowski, Christian Knüpfer, Wolfram Liebermeister, Dagmar Waltemath  
*Cold Spring Harbor Laboratory* (2016-03-21) <https://doi.org/gg939z>  
DOI: [10.1101/044818](https://doi.org/10.1101/044818)
72. **A Coefficient of Agreement for Nominal Scales**  
Jacob Cohen



### 73. Peer review and the publication process

Parveen Azam Ali, Roger Watson

*Nursing Open* (2016-03-16) <https://doi.org/c4g8>

DOI: [10.1002/nop2.51](https://doi.org/10.1002/nop2.51) · PMID: [27708830](https://pubmed.ncbi.nlm.nih.gov/27708830/) · PMCID: [PMC5050543](https://pubmed.ncbi.nlm.nih.gov/PMC5050543/)

### 74. Exploring single-cell data with deep multitasking neural networks

Matthew Amodio, David van Dijk, Krishnan Srinivasan, William S. Chen, Hussein Mohsen, Kevin R. Moon, Allison Campbell, Yujiao Zhao, Xiaomei Wang, Manjunatha Venkataswamy, ... Smita Krishnaswamy

*Nature Methods* (2019-10-07) <https://doi.org/gf9rsg>

DOI: [10.1038/s41592-019-0576-7](https://doi.org/10.1038/s41592-019-0576-7) · PMID: [31591579](https://pubmed.ncbi.nlm.nih.gov/31591579/)

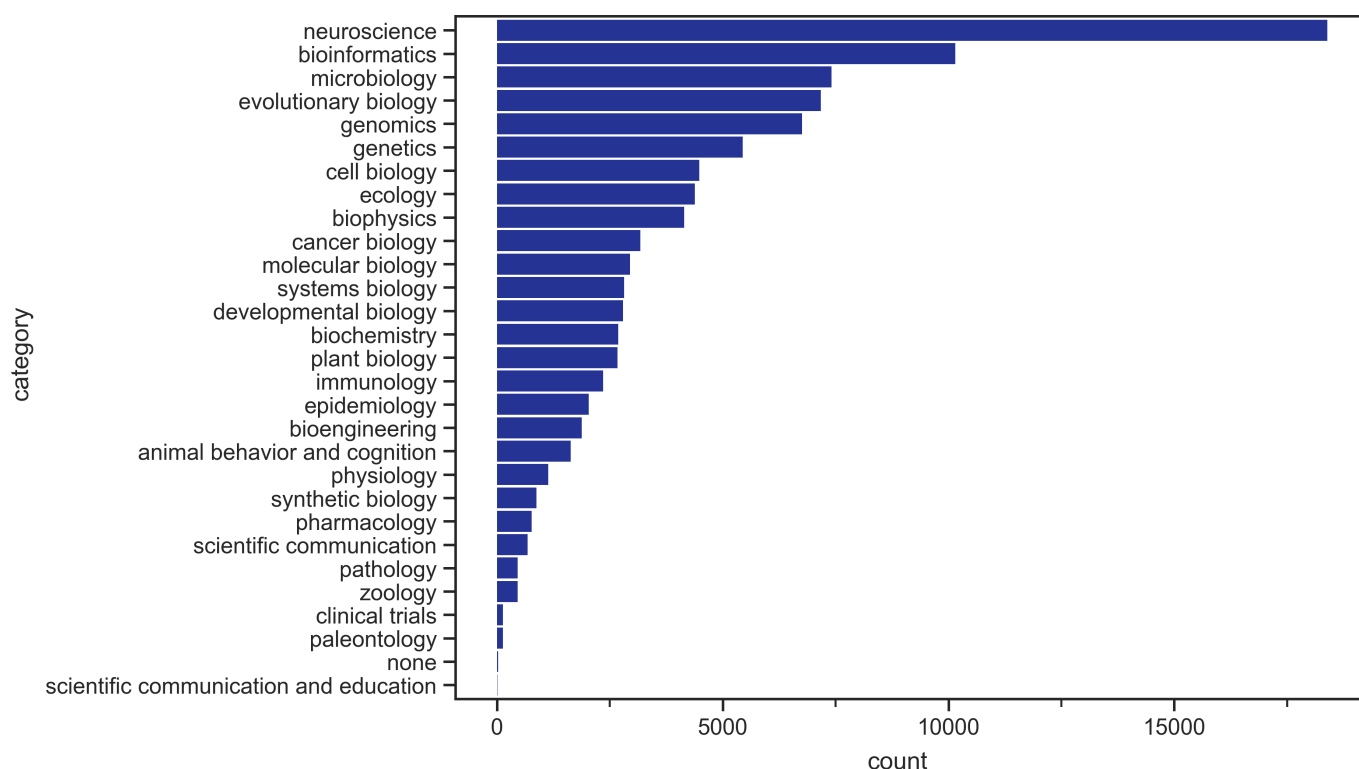
### 75. Quantifying and contextualizing the impact of bioRxiv preprints through automated social media audience segmentation

Jedidiah Carlson, Kelley Harris

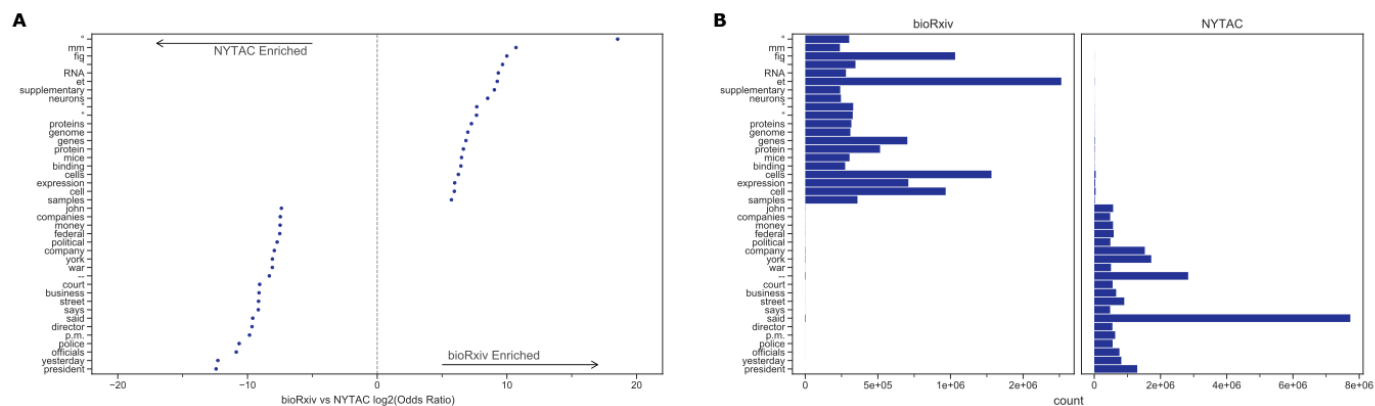
*PLOS Biology* (2020-09-22) <https://doi.org/ghk53x>

DOI: [10.1371/journal.pbio.3000860](https://doi.org/10.1371/journal.pbio.3000860) · PMID: [32960891](https://pubmed.ncbi.nlm.nih.gov/32960891/) · PMCID: [PMC7508356](https://pubmed.ncbi.nlm.nih.gov/PMC7508356/)

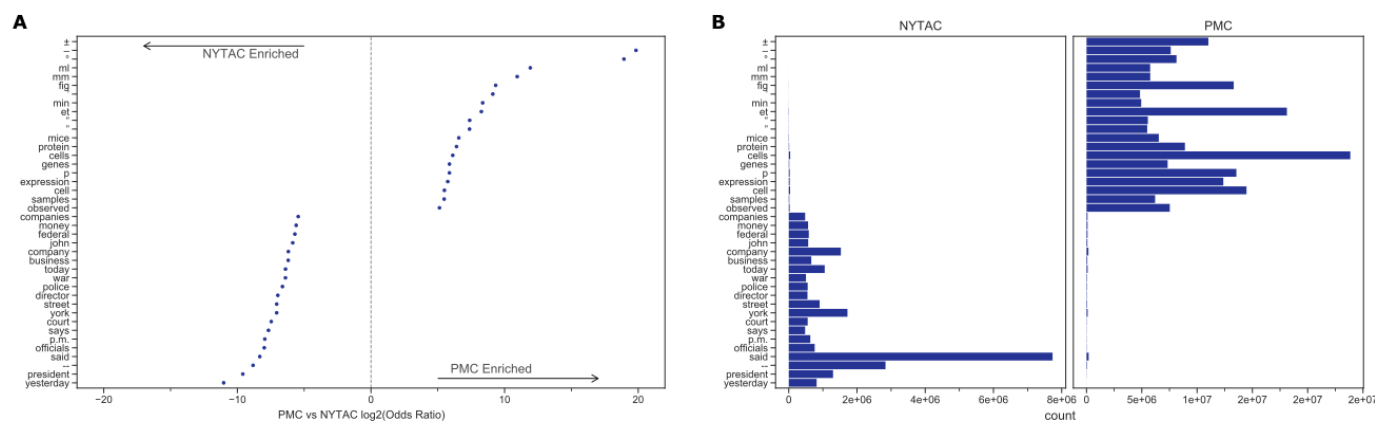
## Supplemental Figures



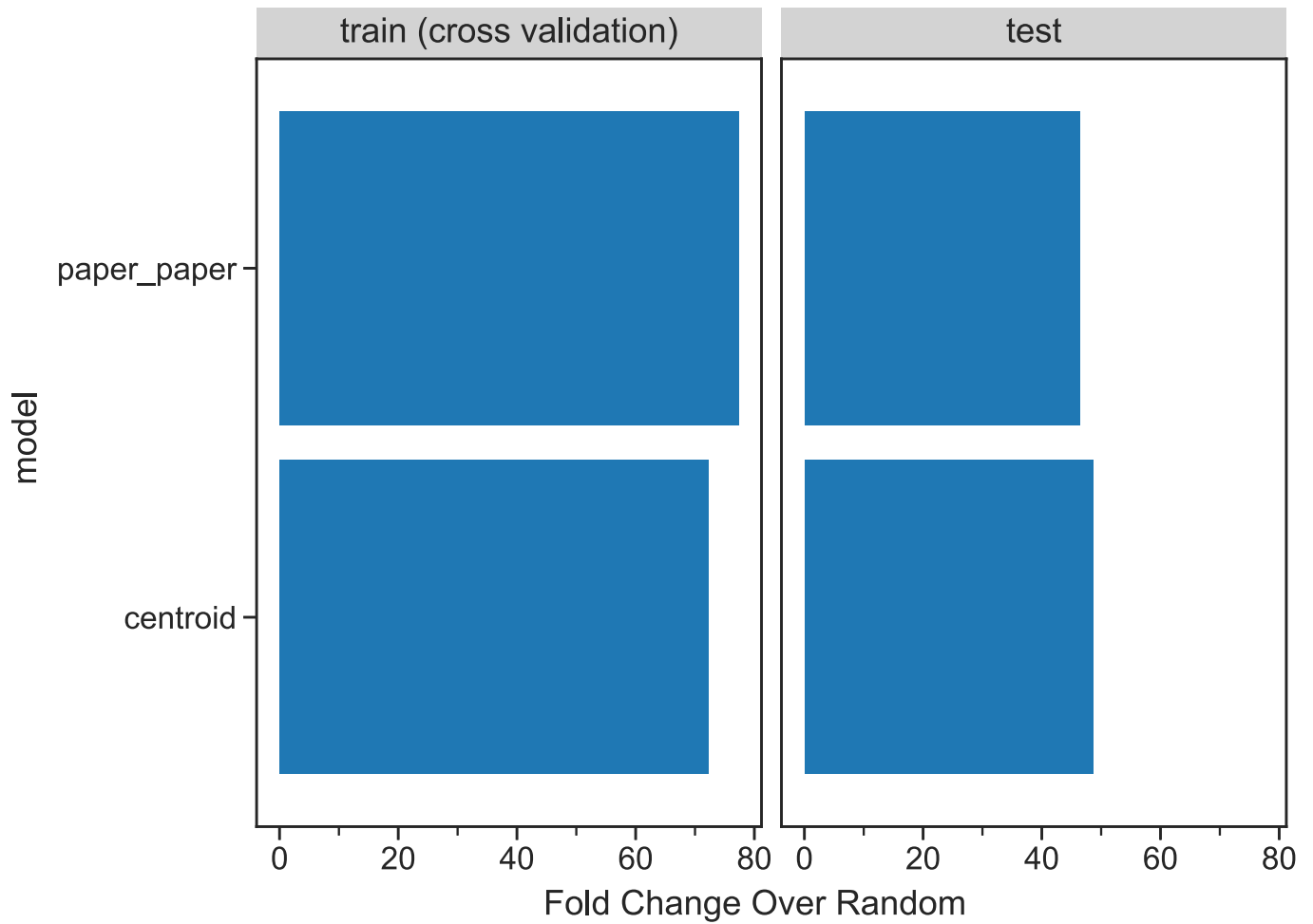
**Figure S1:** Neuroscience and bioinformatics are the two most common author-selected topics for bioRxiv preprints.



**Figure S2:** Topic associated tokens are highly enriched when comparing bioRxiv to the New York Times. The plot on the left (A) is a point range plot of the odds ratio with respect to bioRxiv. Values greater than one indicate a high association with bioRxiv whereas values less than one indicate high association with the New York Times. The dotted line provides a breaking point between both categories. The plot on the right (B) is a bar chart of token frequency appearing in bioRxiv and New York Times respectively.



**Figure S3:** Typesetting symbols and biologically relevant tokens are highly enriched when comparing PubMed Central (PMC) to the New York Times. The plot on the left (A) is a point range plot of the odds ratio with respect to PMC. Values greater than one indicate a high association with PMC whereas values less than one indicate high association with the New York Times. The dotted line provides a breaking point between both categories. The plot on the right (B) is a bar chart of token frequency appearing in PMC and New York Times respectively.



**Figure S4:** Both classifiers outperform the randomized baseline when predicting a paper's journal endpoint. This bargraph shows each model's accuracy in respect to predicting the training and test sets.