

# Linguistic Analysis of the bioRxiv Preprint Landscape

This manuscript ([permalink](#)) was automatically generated from [greenelab/annoxiver manuscript@110749a](#) on October 13, 2020.

## Authors

---

- **David N. Nicholson**

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

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

- **Jane Roe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [janeroe](#)

Department of Something, University of Whatever; Department of Whatever, University of Something

# Abstract

---

## Introduction

---

Preprints are scholarly works that are shared before they have been formally peer reviewed and published. The practice of sharing preprints before peer review has a long history [1]. The longest ongoing use started with physicists in the 1990s [2] and have become an essential tool for the scientific community [3,4,5,6]. Preprints were used in the life sciences community during the 1960s before publisher pressure stopped the practice [1]. Over the past decade preprints have made a resurgence within the life sciences community [7,8]. Common preprint repositories include arXiv [9], bioRxiv [7] and medRxiv [10]; however, there are over 60 different repositories available [11]. Generally, preprints are becoming widely accepted and used within the life sciences community [3,12,13,14], which provides an opportunity to examine their impact.

The scientific community has begun to analyze the impact of preprints in the life sciences. Studies have found that preprints are being posted at a high rate [15]. Studies have also found that preprints are rapidly shared on social media, routinely downloaded, and cited often [16]. Articles with matching preprint versions were cited and discussed more often than articles without them [17,18]. Certain categories of preprints seem to be read and shared differently by both scientists and non-scientists [19]. Across preprint servers, analyses suggest that between two-thirds to three-quarters of preprints are eventually published [8,20]. Time required for a preprint to be published can vary from preprint to preprint; however, studies have found that preprints with a single version often take less time than preprints with multiple versions [21].

Existing work largely focuses on the metadata associated with life sciences preprints. There is still a gap in understanding the language authors use in preprints and how that language shifts during the peer review process. Examining this shift will provide a unique opportunity to ascertain parts of the publishing process and how it impacts the scholarly literature. We sought to understand the language landscape of preprints by performing a linguistic analysis of the *bioRxiv* corpus. We examined textual differences between preprints and published literature by comparing the entire corpus of preprints with articles available in the open access PubMed Central repository. We also examined linguistic differences between preprints and their corresponding published pairs. We found that neural-network derived document embeddings provide a useful space for determining the textual similarity of preprints. In this space a preprint's nearest neighbors are likely to share an eventual publishing venue with the preprint itself. We created a webserver that will display these neighbors for any preprint on *bioRxiv* or *medRxiv*. Close proximity within this space also revealed missing preprint-publication pairs that earlier analyses have been unable to identify. Overall, our linguistic analysis reveals the impact of the publishing process, introduces a method to identify matching preprint-published article pairs, demonstrates that the text content of preprints is related to their eventual publication venue, and provides a more complete picture of the fraction of preprints that are eventually published.

## Methods

---

### Datasets

#### BioRxiv

BioRxiv [7] is a repository for biological and biomedical research preprints. We downloaded an xml snapshot of this repository on February 3, 2020 from bioRxiv's Amazon S3 resource [22] that 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 only the latest one. Preprints in this snapshot were grouped into one of twenty-nine different categories by researchers submitting to *bioRxiv*. Each preprint was also classified either as 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.

## PubMed Central

PubMed Central (PMC) [23] is a repository that contains free-to-read articles. PMC contains two types of contributions: closed access articles from research funded by the United States National Institutes of Health (NIH) appearing after an embargo period and articles published under Gold Open Access [24] publishing schemes. Paper availability within PMC is largely dependent on the journal's participation level [25]. 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 [26], or not participate at all. As of September 2019, PMC had 5,725,819 articles available [27]. Out of these 5 million articles, about 3 million were open access and available for text processing systems [28,29]. We downloaded a snapshot of this open access subset on January 31, 2020. This snapshot contains papers such as literature reviews, book reviews, editorials, case reports, research articles and more; however, we used only research articles.

## Comparing Corpora

We compared bioRxiv against Pubmed Central's Open Access corpus (PMCOA) and the New York Times Annotated corpus (NYTAC) [30] to assess the similarities and differences between bioRxiv, PMCOA and NYTAC. Throughout our analysis we encountered non-word symbols (e.g.,  $\pm$ ), so we refer words and symbols 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 [31] (version 2.2.3) to preprocess all corpora and filtered out 326 spaCy-provided stopwords.

Following the cleaning process, 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 corpus to compare them. We generated a contingency table for each token and calculated the odds ratio from each generated table. We also calculated the 95% confidence interval for each token's odds ratio [32] and measured corpus similarity by calculating the KL divergence across all three corpora.

## Visualizing the Preprint Landscape

### Generate Document Representation

We used gensim [33] (version 3.8.1) to train a word2vec continuous bag of words (CBOW) [34] 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. This document vector is calculated by taking the average of every token present within a given article, ignoring those that were absent from the word2vec model.

### Dimensionality Reduction of Document Embeddings

We used principal component analysis (PCA) [35] to project bioRxiv document vectors into a low dimensional space. We trained this model using scikit-learn's [36] implementation of a randomized

solver [37] with a random seed of 100, output of 50 principal components, and default settings for all other hyperparameters. For each principal component we calculated its cosine similarity with 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

Automated procedures are in place to link preprints to peer reviewed versions and many journals require authors to update preprints with a link to the published version. However, automated procedures at *bioRxiv* are often based on exact matching of certain attributes and authors can forget to establish a link after publication. For example, authors can change the title between a preprint and published version (e.g., [38] and [39]), 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. First, we used CrossRef [40] to identify *bioRxiv* preprints that were linked to a corresponding published article. We filtered out links that contained papers not in PubMed Central's Open Access corpus. Following the preprocessing step, we calculated the distribution of known preprint to published distances by taking 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 candidates to two scientists 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 reviewers. The preprint-publication pairs on which reviewers disagreed were supplied to a third scientist, who carefully reviewed each case and made a final determination. Lastly, 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.

## Journal Recommendation

Determining the best journal venue for a preprint is a non-trivial task as there are too many options for authors to decide. We sought to provide a resource that recommends journals based on a preprint's embedding representation. We illustrate our recommendations as a short list along with a network visualization available at <https://greenelab.github.io/annorxiver-journal-recommender/>. We used a simple k-nearest neighbors approach with Euclidean distance to recommend journals as we sought to examine if embeddings were related to publication venue.

First, we filtered all journals that had fewer than 100 papers in the PMC dataset. A subset of our PMC corpus was directly linked to papers in *bioRxiv* as they had been published as open access articles. We held out this subset and treated it as our gold standard test set. We used the remainder of the PMC corpus for training and initial evaluation via cross validation. We 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 for recommending journals.

The first approach is based on individual paper proximity, which enabled us to provide an example of the specific article or articles that led to the prediction. Conversely, predictions using this technique could be biased due to the overrepresentation of general topic journals. We call this approach the paper-based classifier. This classifier takes a query article that has been projected onto the embedding space constructed by our word2vec model as input and reports the journals of the top ten closest papers. 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.

The second approach is based on close proximity to a journal's centroid. This technique provides recommendations that are more focused on domain-specific publication venues. We call this approach the journal-based classifier. This classifier was trained by computing journal centroids via taking the average embedding of all papers published in each journal. Following the centroid calculation, this classifier takes a query article projected onto the same embedding space as above for input and reports the top ten nearest journals centroids. Both the paper-based classifier and the journal-based classifier were optimized via 10-fold cross validation. We evaluated performance of both classifiers on our gold standard test set of published preprints.

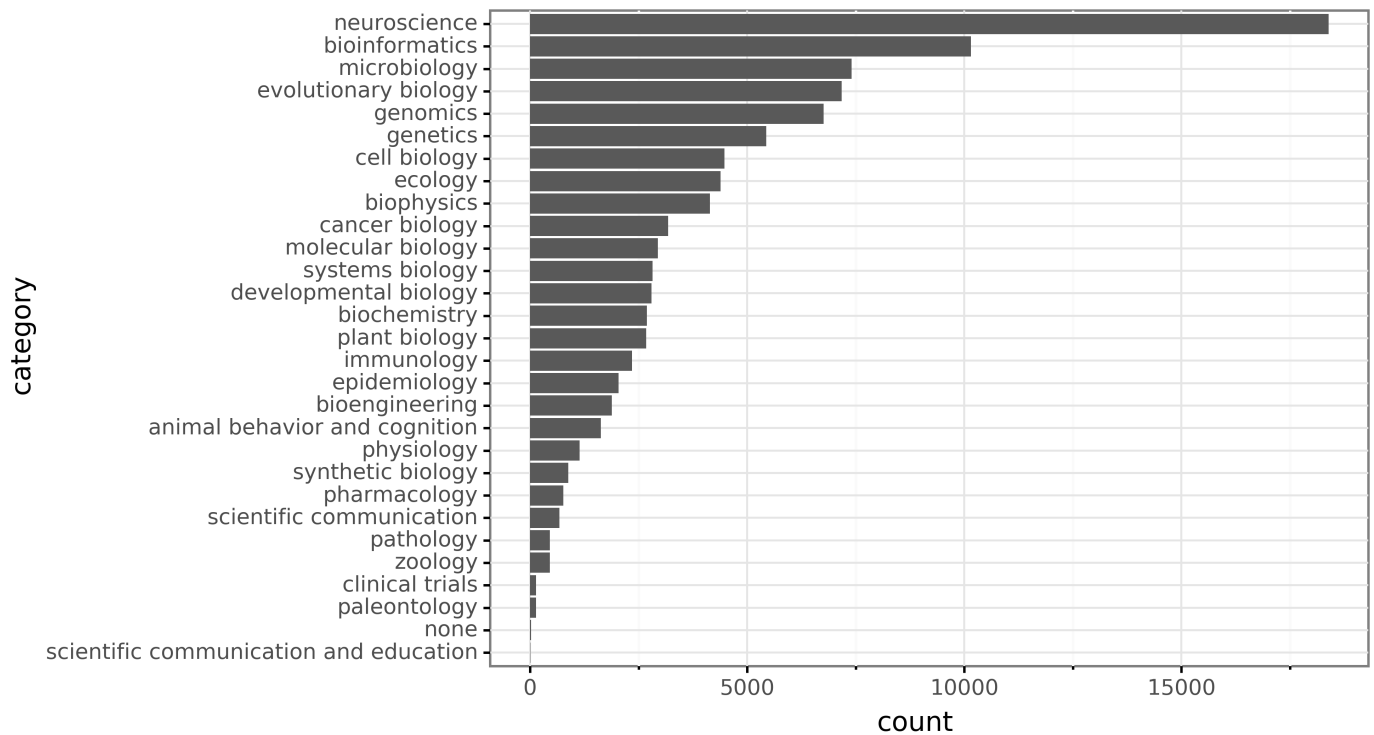
We used SAUCIE [\[41\]](#) to train a model that uses the latent space of a neural network to learn an embedding suitable for visualization. This model enabled us to visualize the PMC corpus and to efficiently embed new papers and preprints within this space. We trained this model using a learning rate of 0.001,  $\lambda_b$  of 0,  $\lambda_c$  of 0.001, and  $\lambda_d$  of 0.001 for 2000 iterations. 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 PubMed Central

#### bioRxiv Repository



**Figure 1:** Neuroscience and bioinformatics are the two most common topics for preprints on bioRxiv. This bar chart depicts the number of preprints that fall into each author-selected topic area.

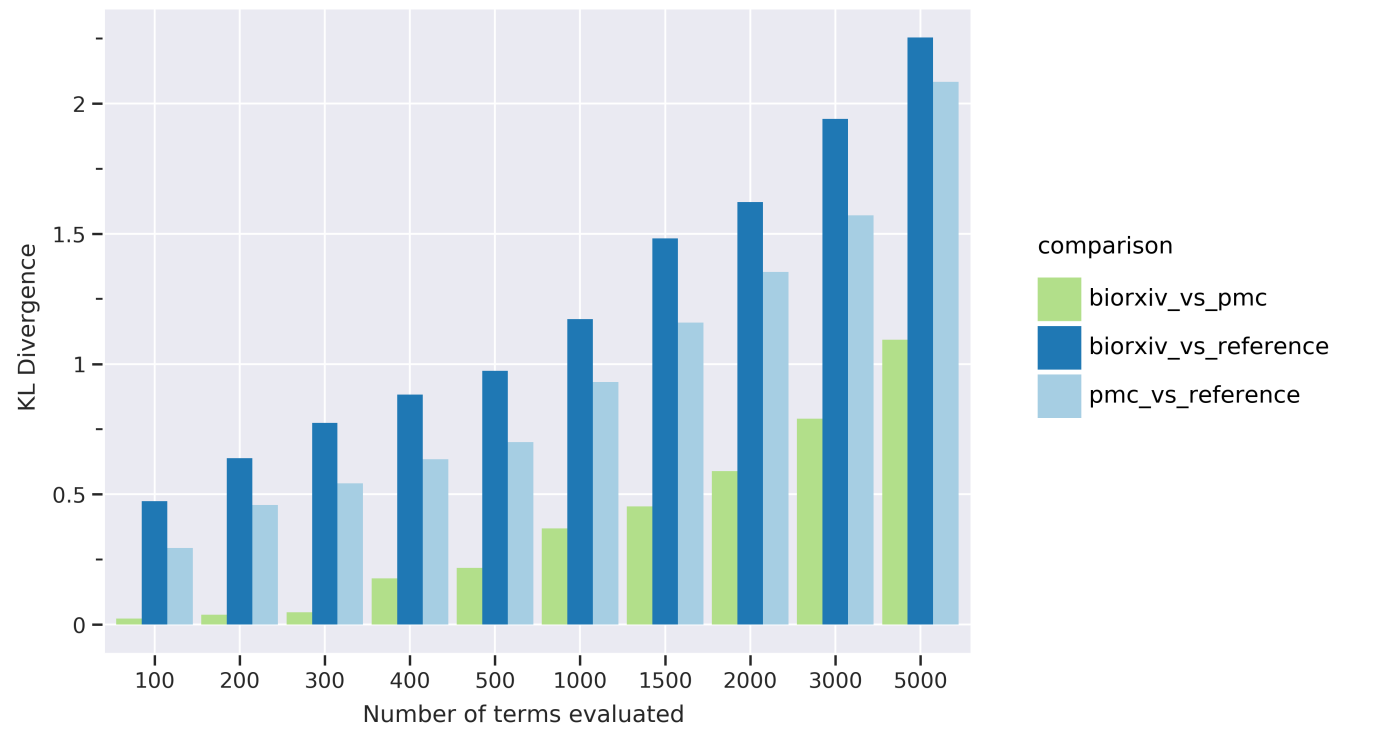
Each preprint on bioRxiv has an author-selected topic area and the predominant area in past reports has been neuroscience [42]. Our analysis of the full text release of bioRxiv confirms this previous finding (Figure 1). The author-selected topic area abundances that we found in the full text largely matched previous studies [42,43]. One exception was microbiology, which has a larger share of preprints than in a previous report from 2018 [42] (Figure 1). 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 were categorized as new results, which is consistent with previous findings [43].

### Global Comparison of bioRxiv and PubMed Central

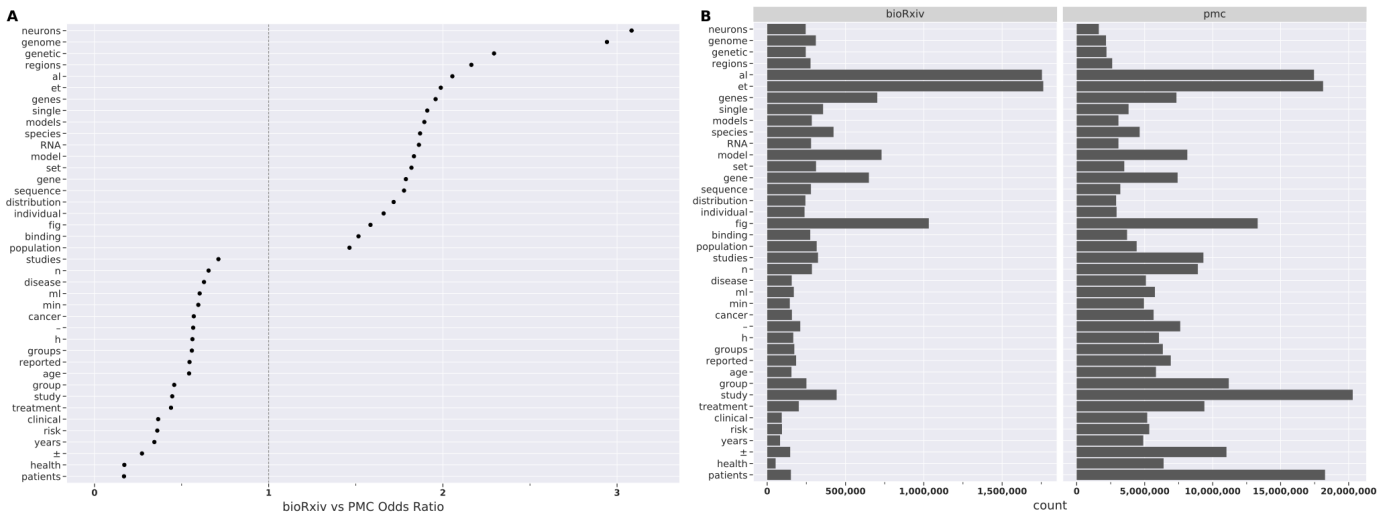
**Table 1:** Generated corpora statistics for all corpus used in this project.

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%

Metric	bioRxiv	PMC	NYTAC
passives	15,012,441	342,407,363	19,472,053
passive%	3.57%	3.98%	1.60%



**Figure 2:** BioRxiv is more similar to PubMed Central than to the reference corpus. This barplot represents the KL divergence between bioRxiv, Pubmed Central and the reference corpus. The y-axis is the KL divergence metric where lower values indicates similar distributions and vice versa for higher values. The x-axis represents the number of highly occurring tokens used to calculate the KL divergence.



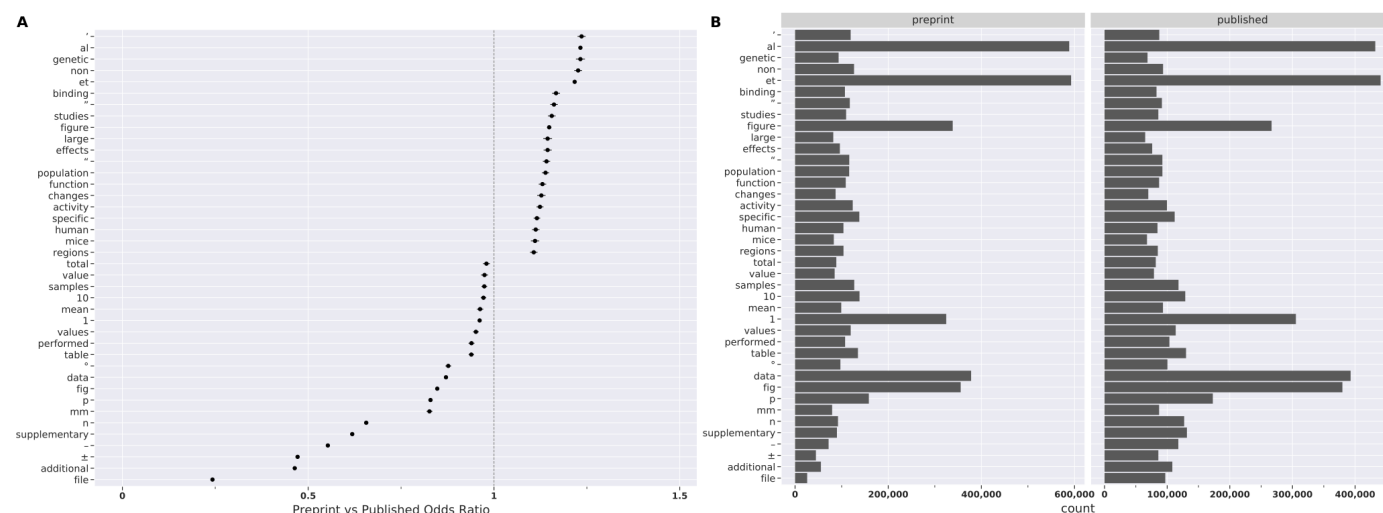
**Figure 3:** BioRxiv is more focused on biological discoveries rather than disease treatments and clinical trials. The plot on the left (A) is a point range plot of the odds ratio with respect to bioRxiv. Values greater than one indicate a high association with bioRxiv whereas values less than one indicate high association with PubMed Central. 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 PMC respectively.

The linguistic style of the bioRxiv corpus differs from the PMC corpus. We compared and contrasted preprints in bioRxiv, published manuscripts in PMC and newspaper articles from the New York Times (NYTAC) against eachother. We refer to NYTAC as our reference corpus for the following analysis. We found that bioRxiv is more similar to PMC than to the reference in terms of token frequencies and



corpora statistics (Figure 2 and Table 1). When comparing bioRxiv and PMC to the reference, topic associated and measurement related tokens appear highly enriched (Supplemental Figures 13 and 14). Furthermore, we found that tokens such as “neuron”, “genome”, “RNA” and “network” had a high odds ratio, while tokens such as “patient”, “health”,  $\pm$ , and “ml” to have a low odds ratio when comparing bioRxiv to PMC (Figure 3). This separation of tokens suggests that articles focused on clinical trials and patient research are more prevalent in PMC than to bioRxiv. This separation also suggests that bioRxiv has a predominance of preprints focused on neuroscience and bioinformatic topics. In regard to writing, citation styles diversify from the familiar “et al.” form as preprints transition through the publication process. Additionally, published articles have an increase of typesetting ( $\pm$ ) and measurement symbols (“ml”, “age”) compared to preprints.

## Published Preprint Differences



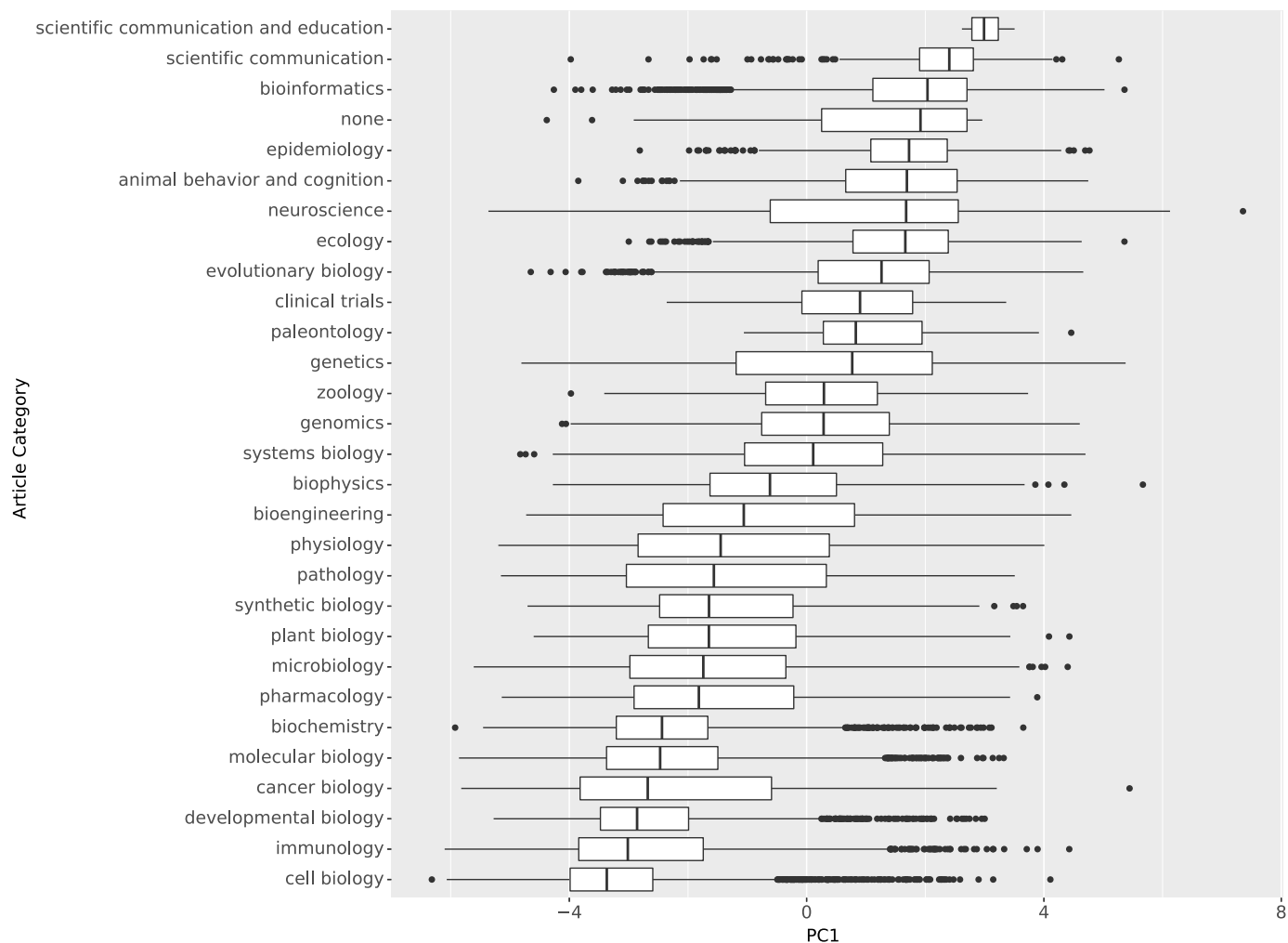
**Figure 4:** Top scoring tokens for preprints are focused on figure citations whereas their published versions are more focused on data availability. The plot on the left (A) is a point range plot of the odds ratio with respect to preprints. Values greater than one indicate a high association with preprints while values less than one indicate a high association with published articles. The dotted line provides a breaking point between both categories. The plot on the right (B) is a barchart of token frequency appearing in preprints and published versions of preprints respectively.

A preprint’s linguistic style can change once a preprint has undergone the revision process prior to being published. We quantified this linguistic difference by calculating the odds ratio of tokens appearing in the union of bioRxiv preprints and their published counterparts within PMC. Tokens with an odds ratio greater than one are mainly centered on paper/figure references and research specific terms (Figure 4). Tokens with an odds ratio of less than one are focused on data availability, and research measurements such as number of cases and controls or significance testing (Figure 4). This enrichment suggests that a key piece in the publication process is verifying that essential parts of research (e.g. data availability, specific measurements) are obvious to future readers within the scientific community.

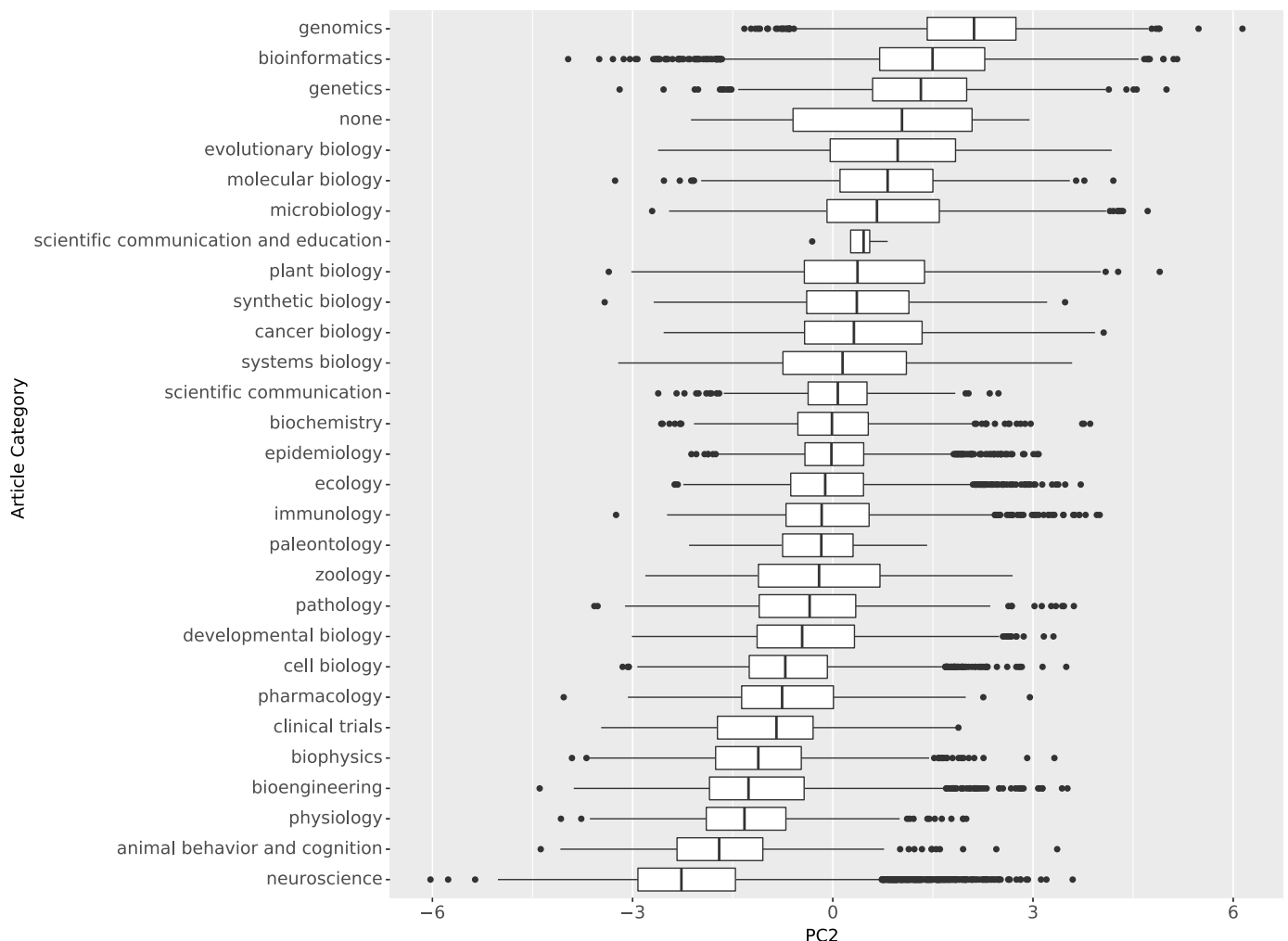
## Topic Analysis of bioRxiv’s Principal Components







**Figure 6:** Preprint categories have a diverse spread of quantitative and molecular biology results. This box plot shows preprints in each article category projected along the PC1 direction. Negative values indicate molecular biology concepts, while positive values indicate quantitative biology concepts.

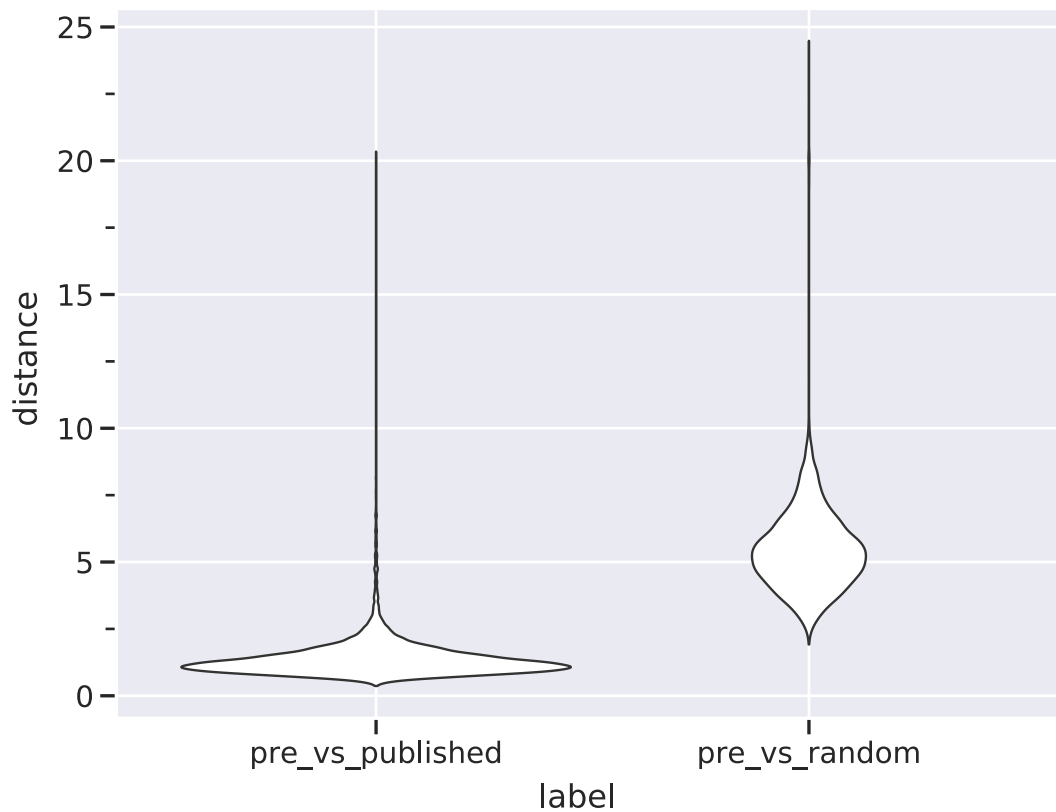


**Figure 7:** The second PC groups neuroscience related preprint categories and bioinformatics related preprint categories together. This box plot shows preprints in each article category projected along the PC2 direction. Negative values indicate neuroscience concepts, while positive values indicate bioinformatic concepts.

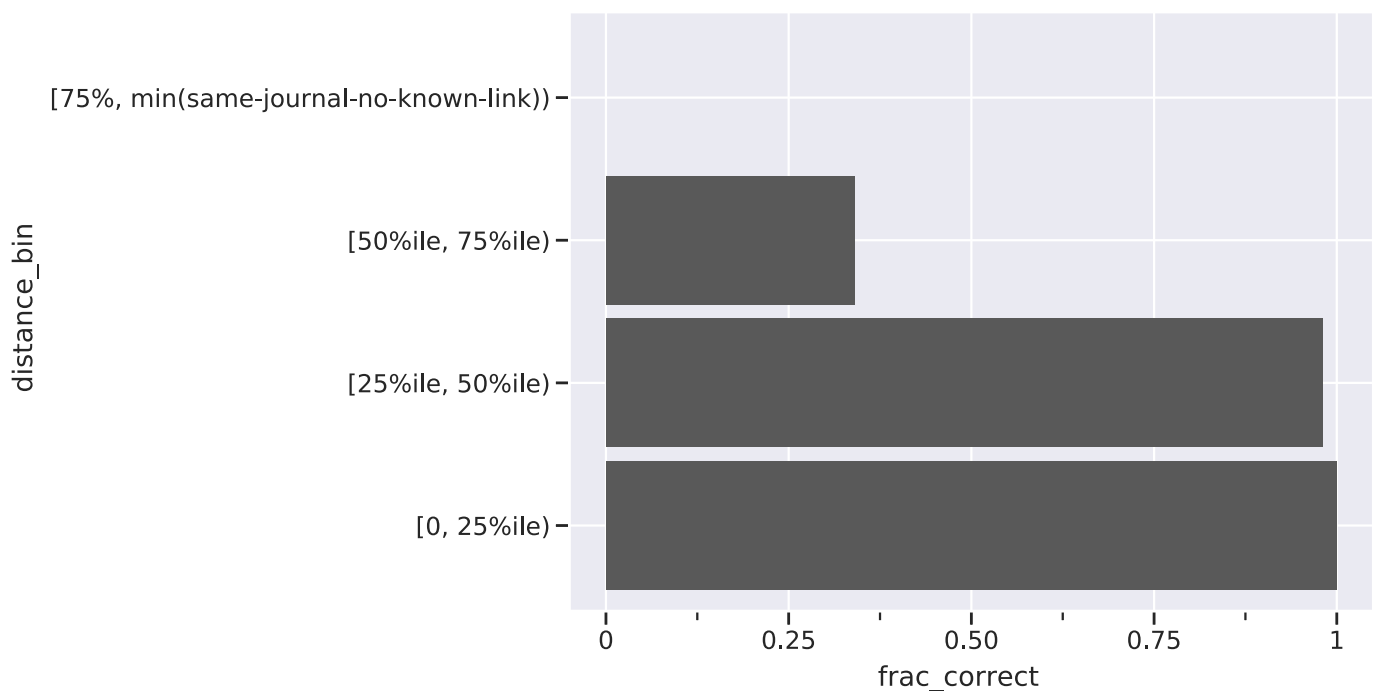
We explored the primary differences between the full text of bioRxiv preprints by performing principal components analysis on generated document embeddings. We visualized the correspondence between tokens and the loadings for each principal component (Figure 5A,C). We also visualized documents projected on selected principal components (Figure 5B). The first principal component separates bioRxiv preprints that encompass molecular biology results with preprints that contain quantitative biology results (Figure 5C). This highlights the bisection of biomedical research where majority of results can be categorized under the molecular biology category or the quantitative biology category. Furthermore, this bisecting trend is evident across individual preprint categories as most categories lie on either side of the first principal component (Figure 6). We also provide example preprints from the systems biology category to reinforce this concept (Supplemental Table 2).

The second principal component represents the concept of neuroscience vs bioinformatics (Figure 5A). This principal component suggests that the bulk of preprints within bioRxiv are largely focused around neuroscience and bioinformatic concepts. This split is evident in Figure 7 as enriched categories along this principal component are quite related to neuroscience (negative end) or bioinformatics (positive end). As with the first principal component we provide example preprints from the systems biology category to reinforce this concept (Supplemental Table 3). More principal component word clouds can be found on our journal recommender website and within our online repository (see Software and Data Availability).

## Identifying preprints that were not linked with their corresponding publications



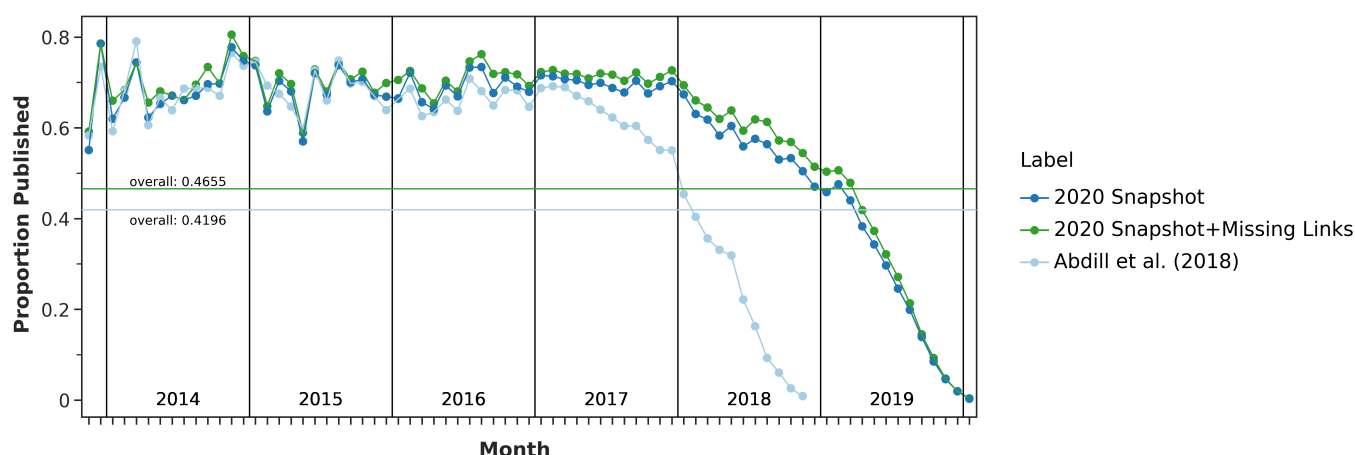
**Figure 8:** The distances between preprints and their published version was on average lower than the distance between preprints and a randomly selected published article in the same journal. This violin plot shows the distribution of distances between both categories.



**Figure 9:** The preprint-published pairs with smaller distances have a high change of being a true match. This bar chart depicts the fraction of true positives over the total number of pairs in each bin. Each bin contains a total of 200 annotated pairs and is based on the percentiles of the preprint-published distribution.

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. Sporadically, there are cases where *bioRxiv* may miss detecting a link or authors may forget to notify *bioRxiv* of their recent publication.

These missing links can make it more difficult to identify the latest version of scientific manuscripts and estimate the fraction of articles that are eventually published [42]. We used distance in the document space to identify preprints without an annotated publication but contained very similar content to published articles. We found that distances between preprints and their corresponding published versions were lower than preprints paired with a random article published in the same journal (Figure 8). This observation suggests that pairs with low embedding distances could be considered a true match, so we separated articles into quantiles based on the distribution of distances between true preprint-publication pairs. We curated 50 potential preprint-publication pairs from each of four quantiles and achieved a high inter-rater reliability of 91.7% (Cohen's Kappa [44]) for this task. Out of these two hundred pairs we found that approximately 98% of pairs with an embedding distance in the 0-25th and 25th-50th percentile bins were true matches (Figure 9). These two bins contained 1,720 preprint-article pairs, suggesting that many preprints have been published but not previously connected with their published versions.



**Figure 10:** The overall fraction of published preprints is higher than originally estimated in [42]. This line plot shows 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. [42]. The dark blue line represents the updated publication rate without missing links added while the dark green line is the updated publication rate with missing links added. The horizontal lines represent the overall proportion of preprints that are published.

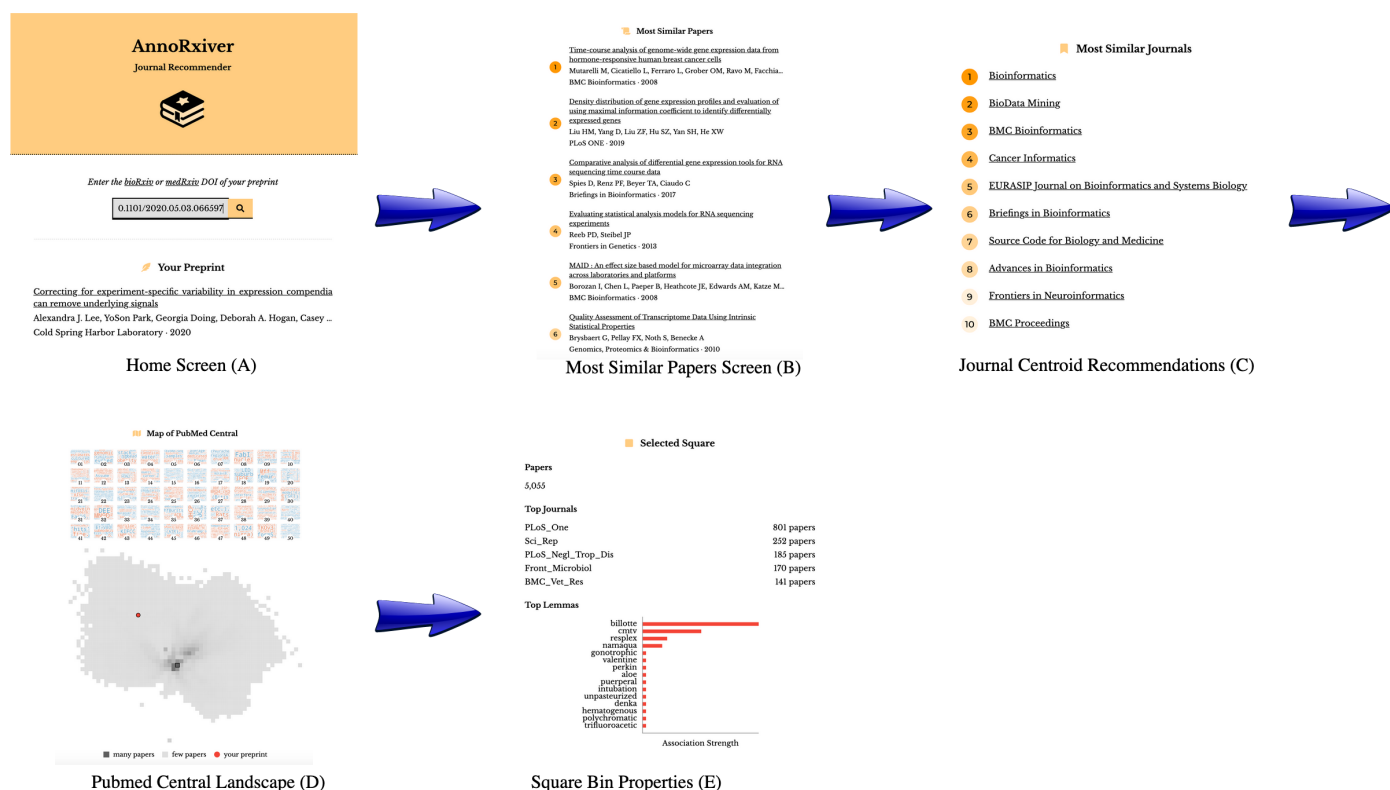
We overlaid these new annotations onto existing annotations to reassess the overall preprint publication rate reported by Abdill et al. [42]. Our filtering criteria were intentionally stringent, so the increased estimate of publication rate amounts to a few percent (Figure 10). Many of these missed annotations were for preprints posted in the 2017-2018 interval. Compared to preprints published in 2019 and later, the preprints posted in 2017-2018 are old enough to have a high chance of being published; however, it is interesting that the rate for older preprints was not observed to be higher.

## Recommending Journals Based on Preprint Representation



**Figure 11:** 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 set.

We sought to identify journals that might publish a preprint based on the text of a paper. We trained two different classifiers to predict the journal endpoints for already published papers. One classifier uses the nearest journal centroids, which attempts to capture the topic area of a journal. The other classifier aims to be more granular and recommends journals based on close proximity of individual papers. Both classifiers achieved a substantial increase over the random baseline; however, our predictors are not perfect (Figure 11). This is expected as our dataset contains 2516 different journals where some journals publish papers that cover very specific topic while others publish papers that have a broad set of covered topics. Our journal centroid classifier performed better than the nearest paper classifier on the held out test set (Figure 11). Overall, our software provides a starting point for authors to use the text of their preprints to identify potentially suitable publication venues.



**Figure 12:** Here is the workflow of the journal recommender web-app. Starting with the homescreen users can paste in a *bioRxiv* or *medRxiv* DOI, which sends a request to biorxiv or medrxiv (A). Next our app preprocesses the preprint and returns a listing of the top ten most similar papers (B) and the top ten closest journals to the query (C). Following the listing, our app manually plots the preprint query onto the Pubmed Central Landscape (D). Lastly, users can click on a square within the landscape, which will show bin statistics as well as associated word-odd ratios (E).

We constructed an online app that provides users with journal suggestions based on their preprint content. Users supply DOIs from *bioRxiv* or *medRxiv*. The application then downloads the article, 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. It also embeds the document into the overall PMC landscape for visualization and allows the user to examine principal components and term enrichment for each bin within the landscape (Figure 12).

## Discussion

We analyzed the language contained used in preprints and examined how it changes through the publication process. We found that *bioRxiv* and PubMed Central (PMC) have similar word frequency distributions, which suggests that the overall manner of writing is consistent with the biomedical literature. At the token level, those most strongly associated with *bioRxiv* are related to neuroscience and bioinformatics, which are also fields that have seen high uptake of preprinting [42]. We noticed that a multitude of preprints highly associated with the first principal component have restrictive or no copyright license (Supplemental Table 2). This finding highlights the ongoing problem of restricted access within the scientific community [45,46]. We also found that the second principal component for our language embedding differentiated neuroscience and bioinformatics papers.

We examined preprints that were textually similar to published articles and found numerous missing links between preprints and their published counterparts. This observation led us to find that the life sciences preprint publication rate is higher than previously estimated (Figure 10). Preprint-publication similarity also predicts journal endpoints with modest performance for already published articles. This observation enabled us to provide a web application that allows users to identify the papers and journals that are most similar to a *bioRxiv* or *medRxiv* preprint.



## Conclusion and Future Directions

---

Our linguistic analysis did not reveal substantial changes in the language during the peer-reviewed publishing process. The tokens most strongly associated with the peer reviewed form, as opposed to the preprint form, were associated with data availability and statistical reporting. We found that embeddings of preprints and publications could be compared and that distance in this space was meaningful in terms of topic area and the journal of eventual publication. Being able to identify similar preprints and publications using text content makes it feasible to begin tackling more detailed questions, and our analytical software is all open source to enable others to build upon them. The analysis of preprints' full text can support new tools that accelerate publishing, integrity checks, and other critically important contributions.

## 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/annorxiver-journal-recommender>. 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>.

## Acknowledgements

---

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*. This work was supported by [Grant GBMF4552](#) from the Gordon Betty Moore Foundation and by NIH T32HG00046, Computational Genomics training grant, from the National Human Genome Research Institute (NHGRI).

# References

---

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

Matthew Cobb

*PLOS Biology* (2017-11-16) <https://doi.org/c6wv>

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/)

**2. Preprint Déjà Vu**

Paul Ginsparg

*The EMBO Journal* (2016-10-19) <https://doi.org/f3r9vf>

DOI: [10.15252/embj.201695531](https://doi.org/10.15252/embj.201695531) · PMID: [27760783](https://pubmed.ncbi.nlm.nih.gov/27760783/) · PMCID: [PMC5167339](https://pubmed.ncbi.nlm.nih.gov/PMC5167339/)

**3. 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/)

**4. 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/)

**5. Preprints for the life sciences**

J. M. Berg, N. Bhalla, P. E. Bourne, M. Chalfie, D. G. Drubin, J. S. Fraser, C. W. Greider, M. Hendricks, C. Jones, R. Kiley, ... C. Wolberger

*Science* (2016-05-19) <https://doi.org/bmp7>

DOI: [10.1126/science.aaf9133](https://doi.org/10.1126/science.aaf9133) · PMID: [27199406](https://pubmed.ncbi.nlm.nih.gov/27199406/)

**6. The rise of preprints in chemistry**

François-Xavier Coudert

*Nature Chemistry* (2020-05-18) <https://doi.org/ghdd64>

DOI: [10.1038/s41557-020-0477-5](https://doi.org/10.1038/s41557-020-0477-5) · PMID: [32424256](https://pubmed.ncbi.nlm.nih.gov/32424256/)

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

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

*bioRxiv* (2019-11-06) <https://doi.org/ggc46z>

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

**8. 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)

**9. 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)

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

11. **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/>
12. **Preprints: An underutilized mechanism to accelerate outbreak science**  
Michael A. Johansson, Nicholas G. Reich, Lauren Ancel Meyers, Marc Lipsitch  
*PLOS Medicine* (2018-04-03) <https://doi.org/gg922h>  
DOI: [10.1371/journal.pmed.1002549](https://doi.org/10.1371/journal.pmed.1002549) · PMID: [29614073](https://pubmed.ncbi.nlm.nih.gov/29614073/) · PMCID: [PMC5882117](https://pubmed.ncbi.nlm.nih.gov/PMC5882117/)
13. **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/)
14. **In praise of preprints**  
Norman K. Fry, Helina Marshall, Tasha Mellins-Cohen  
*Microbial Genomics* (2019-04-01) <https://doi.org/gg3bxc>  
DOI: [10.1099/mgen.0.000259](https://doi.org/10.1099/mgen.0.000259) · PMID: [30938670](https://pubmed.ncbi.nlm.nih.gov/30938670/) · PMCID: [PMC6521583](https://pubmed.ncbi.nlm.nih.gov/PMC6521583/)
15. **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/)
16. **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/)
17. **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)
18. **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/)
19. **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)
20. **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/pra2.175](https://doi.org/10.1002/pra2.175)

**21. 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)

**22. Machine access and text/data mining resources | bioRxiv** <https://www.biorxiv.org/tdm>

**23. 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/pmc/PMC33354/)

**24. 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/pmc/PMC5758455/)

**25. How Papers Get Into PMC** <https://www.ncbi.nlm.nih.gov/pmc/about/submission-methods/>

**26. 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)

**27. PMC Overview** <https://www.ncbi.nlm.nih.gov/pmc/about/intro/>

**28. 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/pmc/PMC6748740/)

**29. 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/pmc/PMC6602571/)

**30. The new york times annotated corpus**

Evan Sandhaus  
*Linguistic Data Consortium, Philadelphia* (2008)

**31. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing**

Matthew Honnibal, Ines Montani  
(2017)

**32. Odds Ratio (OR)**

Steven Tenny, Mary R. Hoffman  
*StatPearls* (2020) <http://www.ncbi.nlm.nih.gov/books/NBK431098/>

**33. Software Framework for Topic Modelling with Large Corpora**

Radim Řehůřek, Petr Sojka

**34. 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>

**35. 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)

**36. 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)

**37. 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>

**38. 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  
*bioRxiv* (2018-07-24) <https://doi.org/gg4hp7>  
DOI: [10.1101/376665](https://doi.org/10.1101/376665)

**39. 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/)

**40. 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)

**41. 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  
*bioRxiv* (2020-03-05) <https://doi.org/gg9353>  
DOI: [10.1101/2020.03.04.975177](https://doi.org/10.1101/2020.03.04.975177)

**42. 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/)

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

Stylianos 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/)

**44. A Coefficient of Agreement for Nominal Scales**

Jacob Cohen

*Educational and Psychological Measurement* (2016-07-02) <https://doi.org/dghsrr>

DOI: [10.1177/001316446002000104](https://doi.org/10.1177/001316446002000104)

**45. Biologists debate how to license preprints**

Lindsay McKenzie

*Nature* (2017-06-16) <https://doi.org/b9fb>

DOI: [10.1038/nature.2017.22161](https://doi.org/10.1038/nature.2017.22161)

**46. The licensing of *bioRxiv* preprints**

Daniel Himmelstein

*Satoshi Village* (2016-12-05) <https://blog.dhimmel.com/biorxiv-licenses/>

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

Fortunato Bianconi, Chiara Antonini, Lorenzo Tomassoni, Paolo Valigi

*bioRxiv* (2017-10-02) <https://doi.org/gg9393>

DOI: [10.1101/197400](https://doi.org/10.1101/197400)

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

Kyemyung Park, Thorsten Prüstel, Yong Lu, John S. Tsang

*bioRxiv* (2019-10-31) <https://doi.org/gg94bm>

DOI: [10.1101/825943](https://doi.org/10.1101/825943)

**49. Notions of similarity for computational biology models**

Ron Henkel, Robert Hoehndorf, Tim Kacprowski, Christian Knüpfer, Wolfram Liebermeister, Dagmar Waltemath

*bioRxiv* (2016-03-21) <https://doi.org/gg939z>

DOI: [10.1101/044818](https://doi.org/10.1101/044818)

**50. 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

*bioRxiv* (2019-09-18) <https://doi.org/gg94bj>

DOI: [10.1101/769299](https://doi.org/10.1101/769299)

**51. SBpipe: a collection of pipelines for automating repetitive simulation and analysis tasks**

Piero Dalle Pezze, Nicolas Le Novère

*bioRxiv* (2017-02-09) <https://doi.org/gg9392>

DOI: [10.1101/107250](https://doi.org/10.1101/107250)

**52. Spatiotemporal proteomics uncovers cathepsin-dependent host cell death during bacterial infection**

Joel Selkrig, Nan Li, Jacob Bobonis, Annika Hausmann, Anna Sueki, Haruna Imamura, Bachir El Debs, Gianluca Sigismondo, Bogdan I. Florea, Herman S. Overkleeft, ... Athanasios Typas

*bioRxiv* (2018-11-07) <https://doi.org/gg94bc>  
DOI: [10.1101/455048](https://doi.org/10.1101/455048)

**53. 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

*bioRxiv* (2018-07-19) <https://doi.org/gg9398>  
DOI: [10.1101/371922](https://doi.org/10.1101/371922)

**54. 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

*bioRxiv* (2019-08-13) <https://doi.org/gg94bf>  
DOI: [10.1101/733162](https://doi.org/10.1101/733162)

**55. 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'Riordan, Lynda Zeboudj, Nick Guisot, Peter Bunyard, David R. Greaves, Christoph Thiemermann

*bioRxiv* (2019-08-28) <https://doi.org/gg94bg>  
DOI: [10.1101/745943](https://doi.org/10.1101/745943)

**56. 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

*bioRxiv* (2019-09-01) <https://doi.org/gg94bh>  
DOI: [10.1101/754572](https://doi.org/10.1101/754572)

**57. Pangenome Analysis of Enterobacteria Reveals Richness of Secondary Metabolite Gene Clusters and their Associated Gene Sets**

Omkar S. Mohite, Colton J. Lloyd, Jonathan M. Monk, Tilmann Weber, Bernhard O. Palsson

*bioRxiv* (2019-09-25) <https://doi.org/gg94bk>  
DOI: [10.1101/781328](https://doi.org/10.1101/781328)

**58. QTG-Finder: a machine-learning based algorithm to prioritize causal genes of quantitative trait loci**

Fan Lin, Jue Fan, Seung Y. Rhee

*bioRxiv* (2019-04-29) <https://doi.org/gg94bd>  
DOI: [10.1101/484204](https://doi.org/10.1101/484204)

**59. Identification of candidate genes underlying nodulation-specific phenotypes in *Medicago truncatula* through integration of genome-wide association studies and co-expression networks**

Jean-Michel Michno, Liana T. Burghardt, Junqi Liu, Joseph R. Jeffers, Peter Tiffin, Robert M. Stupar, Chad L. Myers

*bioRxiv* (2018-08-16) <https://doi.org/gg94bb>  
DOI: [10.1101/392779](https://doi.org/10.1101/392779)

**60. Raw sequence to target gene prediction: An integrated inference pipeline for ChIP-seq and RNA-seq datasets**

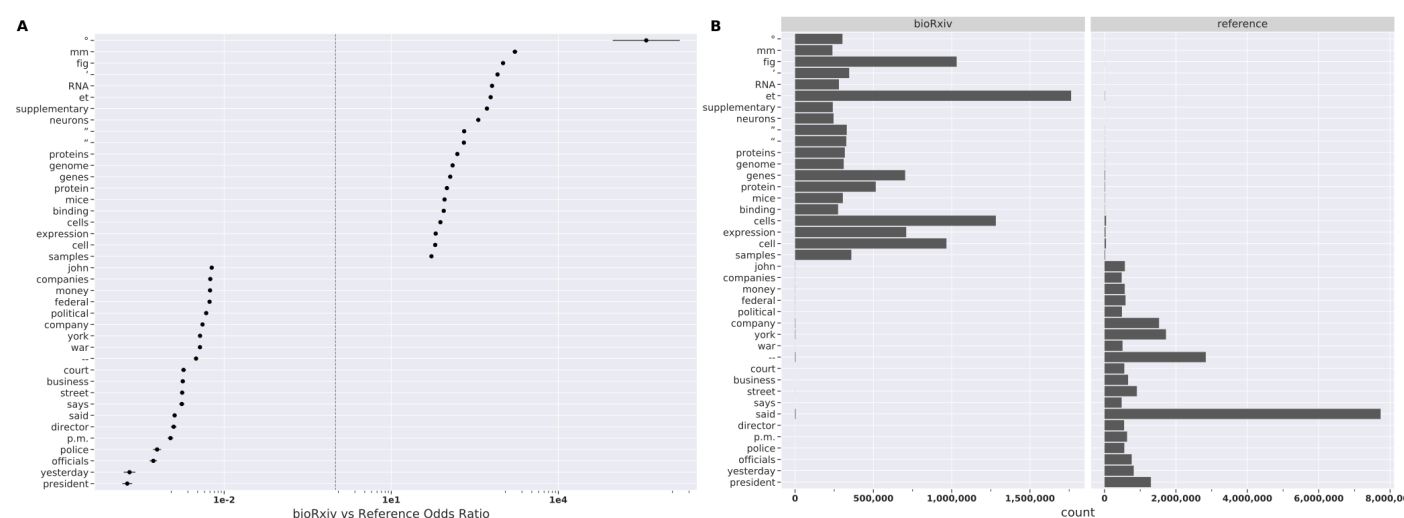
Nisar Wani, Khalid Raza

*bioRxiv* (2017-11-16) <https://doi.org/gg9394>  
DOI: [10.1101/220152](https://doi.org/10.1101/220152)



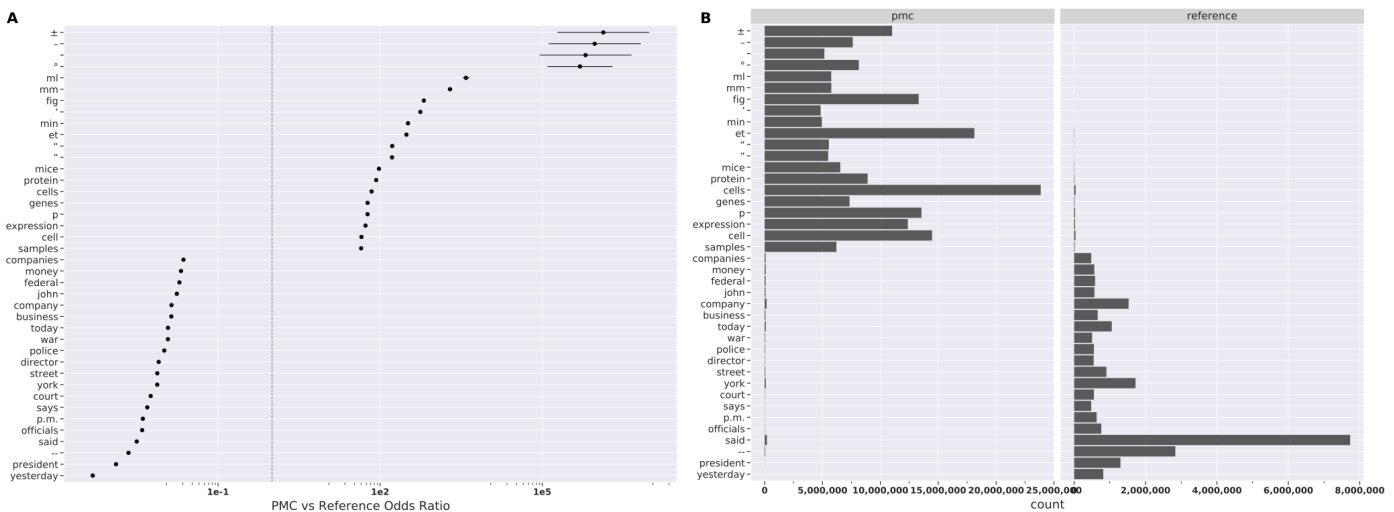
61. **The y-ome defines the thirty-four percent of *Escherichia coli* genes that lack experimental evidence of function**  
Sankha Ghatak, Zachary A. King, Anand Sastry, Bernhard O. Palsson  
*bioRxiv* (2018-12-03) <https://doi.org/gg9396>  
DOI: [10.1101/328591](https://doi.org/10.1101/328591)
62. **The effects of time-varying temperature on delays in genetic networks**  
Marcella M Gomez, Richard M Murray, Matthew R Bennett  
*bioRxiv* (2015-09-24) <https://doi.org/gg939x>  
DOI: [10.1101/019687](https://doi.org/10.1101/019687)
63. **An analog to digital converter creates nuclear localization pulses in yeast calcium signaling**  
Ian S Hsu, Bob Strome, Sergey Plotnikov, Alan M Moses  
*bioRxiv* (2018-06-28) <https://doi.org/gg9397>  
DOI: [10.1101/357939](https://doi.org/10.1101/357939)
64. **Nicotinic modulation of hierarchal inhibitory control over prefrontal cortex resting state dynamics: modeling of genetic modification and schizophreniarelated pathology**  
Marie Rooy, Fani Koukoulis, Uwe Maskos, Boris Gutkin  
*bioRxiv* (2018-04-13) <https://doi.org/gg9395>  
DOI: [10.1101/301051](https://doi.org/10.1101/301051)
65. **Electrical propagation of vasodilatory signals in capillary networks**  
Pilhwa Lee  
*bioRxiv* (2019-11-13) <https://doi.org/gg94bn>  
DOI: [10.1101/840280](https://doi.org/10.1101/840280)
66. **Dendritic spine geometry and spine apparatus organization govern the spatiotemporal dynamics of calcium**  
Miriam Bell, Tom Bartol, Terrence Sejnowski, Padmini Rangamani  
*bioRxiv* (2019-05-29) <https://doi.org/gg9399>  
DOI: [10.1101/386367](https://doi.org/10.1101/386367)

## Supplemental Figures



**Figure 13:** 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.

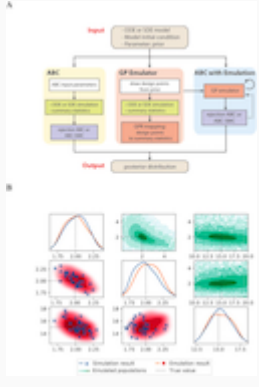
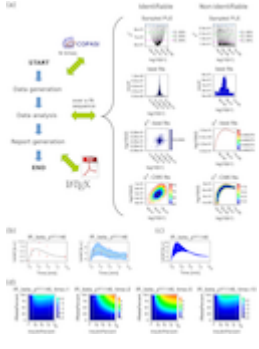
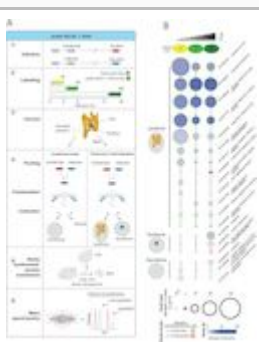
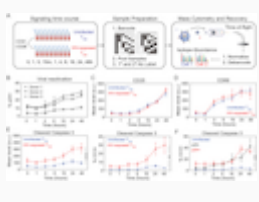
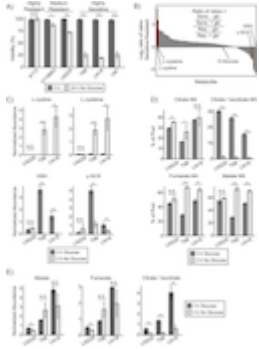


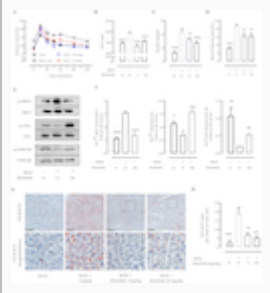
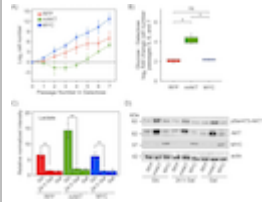
**Figure 14:** Typesetting symbols and biologically relevant tokens are highly enriched when comparing PubMed Central (PMC) to the New York Times. The plot on the left (A) is a point range plot of the odds ratio with respect to PMC. Values greater than one indicate a high association with PMC whereas 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.

## Supplemental Tables

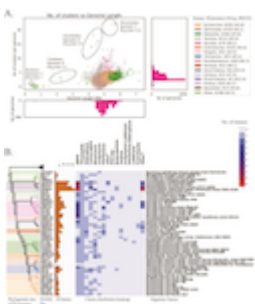
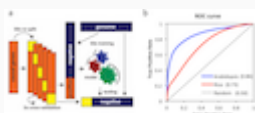
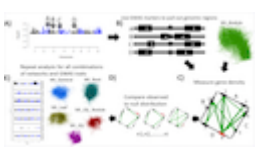
**Table 2:** Top and bottom five systems biology preprints projected onto the PC1 direction. These preprints contain quantitative and molecular biology concepts respectively.

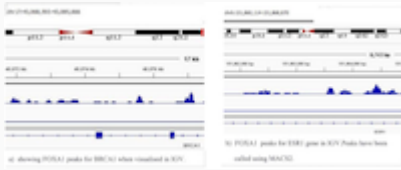
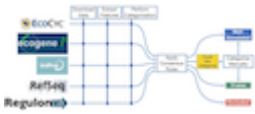

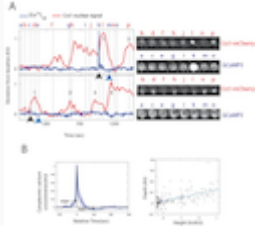
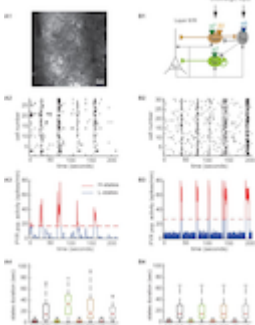
Title [citation]	PC_1	License	Figure Thumbnail	Figure Link
Conditional Robust Calibration (CRC): a new computational Bayesian methodology for model parameters estimation and identifiability analysis [47]	4.700554908074704	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2017/10/02/197400/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2017/10/02/197400/F1.large.jpg</a>
Machine learning of stochastic gene network phenotypes [48]	4.410660604449826	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/10/31/825943/F5.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/10/31/825943/F5.large.jpg</a>
Notions of similarity for computational biology models [49]	4.355295926618207	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2016/03/21/044818/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2016/03/21/044818/F1.large.jpg</a>

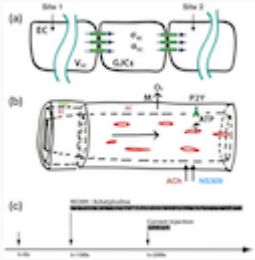
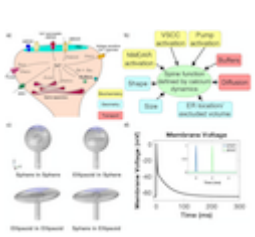
Title [citation]	PC_1	License	Figure Thumbnail	Figure Link
GpABC: a Julia package for approximate Bayesian computation with Gaussian process emulation [50]	4.351517618262304	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/09/18/769299/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/09/18/769299/F1.large.jpg</a>
SBpipe: a collection of pipelines for automating repetitive simulation and analysis tasks [51]	4.321847854182741	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2017/02/09/107250/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2017/02/09/107250/F1.large.jpg</a>
Spatiotemporal proteomics uncovers cathepsin-dependent host cell death during bacterial infection [52]	-4.263964235099807	CC-BY-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/11/07/455048/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/11/07/455048/F1.large.jpg</a>
Systems analysis by mass cytometry identifies susceptibility of latent HIV-infected T cells to targeting of p38 and mTOR pathways [53]	-4.279016673409032	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/07/19/371922/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/07/19/371922/F1.large.jpg</a>
NADPH consumption by L-cystine reduction creates a metabolic vulnerability upon glucose deprivation [54]	-4.592209988884236	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/08/13/733162/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/08/13/733162/F1.large.jpg</a>

Title [citation]	PC_1	License	Figure Thumbnail	Figure Link
Inhibition of Bruton's tyrosine kinase reduces NF- $\kappa$ B and NLRP3 inflammasome activity preventing insulin resistance and microvascular disease [55]	-4.736613689905791	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/08/28/745943/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/08/28/745943/F1.large.jpg</a>
AKT but not MYC promotes reactive oxygen species-mediated cell death in oxidative culture [56]	-4.826793742506695	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/09/01/754572/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/09/01/754572/F1.large.jpg</a>

**Table 3:** Top and bottom five systems biology preprints projected onto the PC2 direction. These preprints contain bioinformatics and neuroscience concepts respectively.

Title [citation]	PC_2	License	Figure Thumbnail	Figure Link
Pangenome Analysis of Enterobacteria Reveals Richness of Secondary Metabolite Gene Clusters and their Associated Gene Sets [57]	3.586570265943883	CC-BY-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/09/25/781328/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/09/25/781328/F1.large.jpg</a>
QTG-Finder: a machine-learning based algorithm to prioritize causal genes of quantitative trait loci [58]	3.470388383023157	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/04/29/484204/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/04/29/484204/F1.large.jpg</a>
Identification of candidate genes underlying nodulation-specific phenotypes in Medicago truncatula through integration of genome-wide association studies and co-expression networks [59]	3.3814906334073953	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/08/16/392779/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/08/16/392779/F1.large.jpg</a>

Title [citation]	PC_2	License	Figure Thumbnail	Figure Link
Raw sequence to target gene prediction: An integrated inference pipeline for ChIP-seq and RNA-seq datasets [60]	3.3632576028389742	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2017/11/16/220152/F3.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2017/11/16/220152/F3.large.jpg</a>
The y-ome defines the thirty-four percent of Escherichia coli genes that lack experimental evidence of function [61]	3.28742786641417	CC-BY		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/12/03/328591/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/12/03/328591/F1.large.jpg</a>
The effects of time-varying temperature on delays in genetic networks [62]	-2.7047102478958056	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2015/09/24/019687/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2015/09/24/019687/F1.large.jpg</a>
An analog to digital converter creates nuclear localization pulses in yeast calcium signaling [63]	-2.775745000260895	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/06/28/357939/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/06/28/357939/F1.large.jpg</a>
Nicotinic modulation of hierarchal inhibitory control over prefrontal cortex resting state dynamics: modeling of genetic modification and schizophreniarelated pathology [64]	-3.047342382798414	None		<a href="https://www.biorxiv.org/content/biorxiv/early/2018/04/13/301051/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2018/04/13/301051/F1.large.jpg</a>

Title [citation]	PC_2	License	Figure Thumbnail	Figure Link
Electrical propagation of vasodilatory signals in capillary networks [65]	-3.107715578793087	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/11/13/840280/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/11/13/840280/F1.large.jpg</a>
Dendritic spine geometry and spine apparatus organization govern the spatiotemporal dynamics of calcium [66]	-3.21533499072831	CC-BY-NC-ND		<a href="https://www.biorxiv.org/content/biorxiv/early/2019/05/29/386367/F1.large.jpg">https://www.biorxiv.org/content/biorxiv/early/2019/05/29/386367/F1.large.jpg</a>