

Analysis of science journalism reveals gender and regional disparities in coverage

This manuscript ([permalink](#)) was automatically generated from greenelab/nature_news_manuscript@af52066 on March 16, 2022.

Authors

- **Natalie R. Davidson**

 [0000-0002-1745-8072](#) ·  [nrosed](#) ·  [n_rose_d](#)

University of Colorado School of Medicine, Aurora, Colorado, United States of America; Center for Health AI · Funded by The Gordon and Betty Moore Foundation (GBMF 4552)

- **Casey S. Greene**

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

University of Colorado School of Medicine, Aurora, Colorado, United States of America; Center for Health AI · Funded by The Gordon and Betty Moore Foundation (GBMF 4552)

Abstract

Science journalism is a critical way in which the public can remain informed and benefit from new scientific findings. Such journalism also shapes the public's view of the current state of scientific findings and legitimizes experts. Those covering science can only cite and quote a limited number of sources. Sources may be identified by the journalist's research or by recommendations by other scientists. In both cases, biases may influence who is identified and ultimately included as an expert. To analyze possible biases in science journalism, we analyzed 22,001 non-research articles published by *Nature*. We chose to analyze *Nature* non-research articles since its research articles provide a natural comparator. Our analysis considered two possible sources of disparity: gender and name origin. To explore these sources of disparity, we extracted cited authors' names as well as extracted names of quoted speakers. While citations and quotations within a piece do not reflect the entire information-gathering process, they can provide insight into the demographics of visible sources. We then used the extracted names to predict gender and name origin of the cited authors and speakers.

In order to appropriately quantify the level of difference, we must identify a suitable reference set for comparison. We chose first and last authors within primary research articles in *Nature* and a subset of *Springer Nature* articles in the same time period as our comparator. In our analysis, we found a skew towards male quotation in *Nature* science journalism-related articles. However, quotation is trending toward equal representation at a faster rate than first and last authorship in academic publishing. Interestingly, we found that the gender disparity in *Nature* quotes was column-dependent, with the "Career Features" column reaching gender parity. Our name origin analysis found a significant over-representation of names with predicted Celtic/English origin and under-representation of names with a predicted East Asian origin. This finding was observed both in extracted quotes and journal citations, but dampened in citations.

This manuscript is a template (aka "rootstock") for [Manubot](#), a tool for writing scholarly manuscripts. Use this template as a starting point for your manuscript.

The rest of this document is a full list of formatting elements/features supported by Manubot. Compare the input (`.md` files in the `/content` directory) to the output you see below.

Basic formatting

Bold text

Semi-bold text

Centered text

Right-aligned text

Italic text

Combined *italics* and **bold**

~~Strikethrough~~

1. Ordered list item
2. Ordered list item

- a. Sub-item
 - b. Sub-item
 - i. Sub-sub-item
 - 3. Ordered list item
 - a. Sub-item
- List item
 - List item
 - List item

subscript: H₂O is a liquid

superscript: 2¹⁰ is 1024.

[unicode superscripts](#)⁰¹²³⁴⁵⁶⁷⁸⁹

[unicode subscripts](#)₀₁₂₃₄₅₆₇₈₉

A long paragraph of text. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Putting each sentence on its own line has numerous benefits with regard to [editing](#) and [version control](#).

Line break without starting a new paragraph by putting two spaces at end of line.

Document organization

Document section headings:

Heading 1

Heading 2

Heading 3

Heading 4

Heading 5

Heading 6

A heading centered on its own printed page

Horizontal rule:

Heading 1's are recommended to be reserved for the title of the manuscript.

Heading 2's are recommended for broad sections such as *Abstract*, *Methods*, *Conclusion*, etc.

Heading 3's and Heading 4's are recommended for sub-sections.

Links

Bare URL link: <https://manubot.org>

[Long link with lots of words and stuff and junk and bleep and blah and stuff and other stuff and more stuff yeah](#)

[Link with text](#)

[Link with hover text](#)

[Link by reference](#)

Citations

Citation by DOI [1].

Citation by PubMed Central ID [2].

Citation by PubMed ID [3].

Citation by Wikidata ID [4].

Citation by ISBN [5].

Citation by URL [6].

Citation by alias [7].

Multiple citations can be put inside the same set of brackets [1,5,7]. Manubot plugins provide easier, more convenient visualization of and navigation between citations [2,3,7,8].

Citation tags (i.e. aliases) can be defined in their own paragraphs using Markdown's reference link syntax:

Referencing figures, tables, equations

Figure 1

Figure 2

[Figure 3](#)

[Figure 4](#)

[Table 1](#)

[Equation 1](#)

[Equation 2](#)

Quotes and code

Quoted text

Quoted block of text

Two roads diverged in a wood, and I—
I took the one less traveled by,
And that has made all the difference.

Code `in the middle` of normal text, aka `inline code`.

Code block with Python syntax highlighting:

```
from manubot.cite.doi import expand_short_doi

def test_expand_short_doi():
    doi = expand_short_doi("10/c3bp")
    # a string too long to fit within page:
    assert doi == "10.25313/2524-2695-2018-3-vliyanie-enhansera-copia-i-
        insulyatora-gypsy-na-sintez-ernk-modifikatsii-hromatina-i-
        svyazyvanie-insulyatornyh-belkov-vtransfetsirovannyh-geneticheskikh-
        konstruktsiyah"
```

Code block with no syntax highlighting:

```
Exporting HTML manuscript
Exporting DOCX manuscript
Exporting PDF manuscript
```

Figures

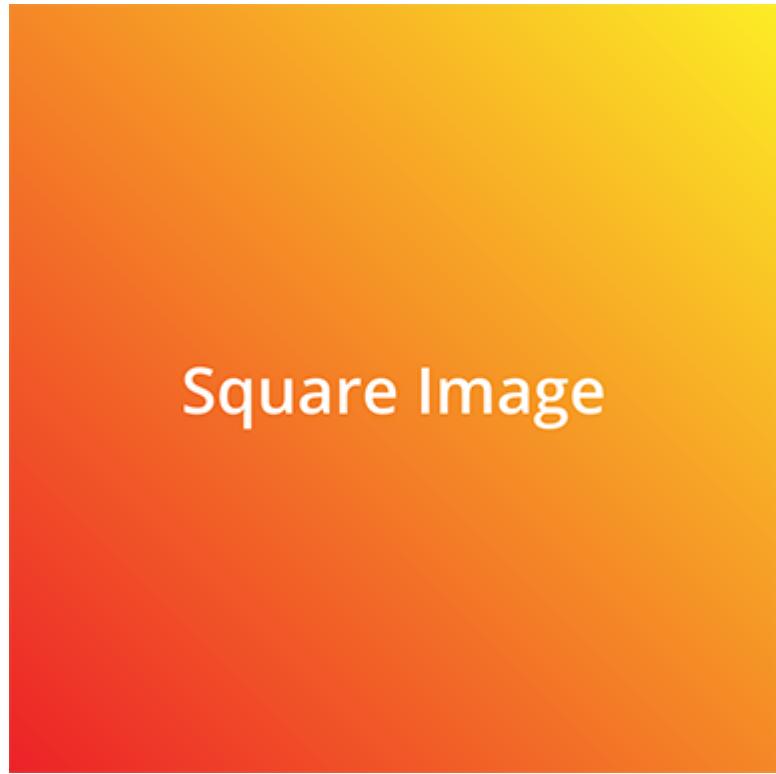


Figure 1: A square image at actual size and with a bottom caption. Loaded from the latest version of image on GitHub.

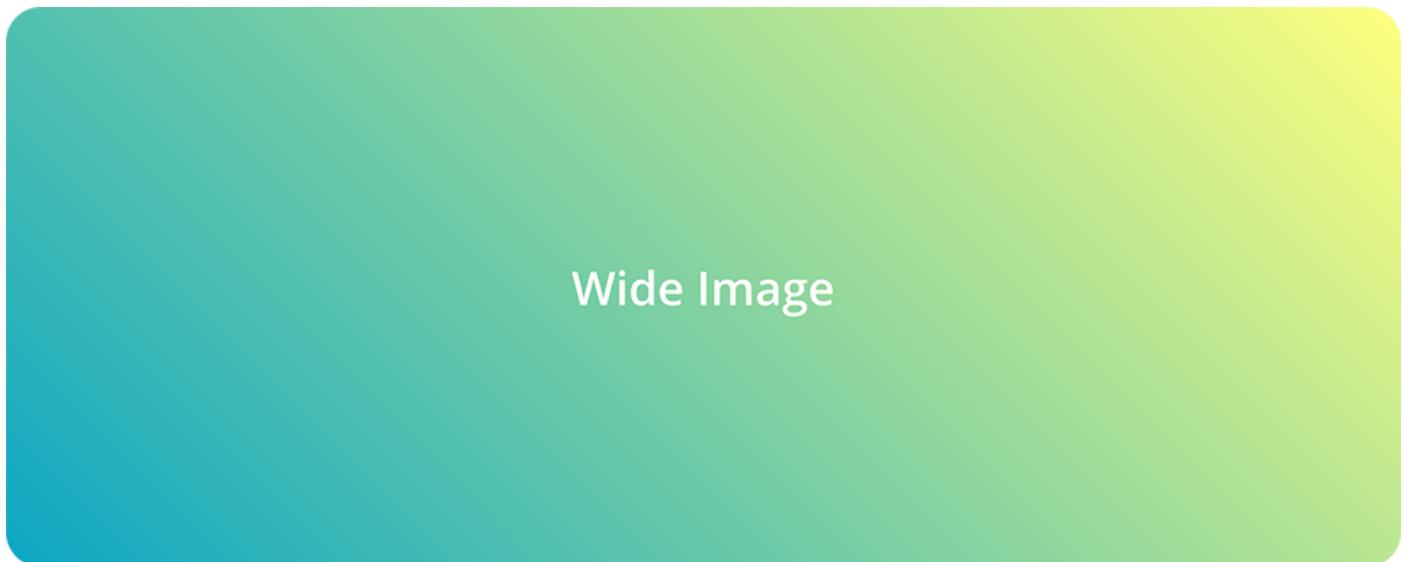


Figure 2: An image too wide to fit within page at full size. Loaded from a specific (hashed) version of the image on GitHub.

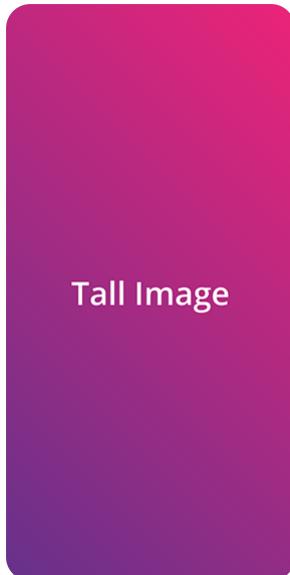


Figure 3: A tall image with a specified height. Loaded from a specific (hashed) version of the image on GitHub.

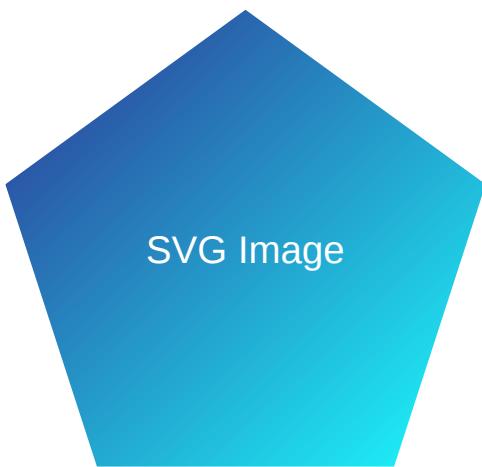


Figure 4: A vector .svg image loaded from GitHub. The parameter `sanitize=true` is necessary to properly load SVGs hosted via GitHub URLs. White background specified to serve as a backdrop for transparent sections of the image.

Tables

Table 1: A table with a top caption and specified relative column widths.

Bowling Scores	Jane	John	Alice	Bob
Game 1	150	187	210	105
Game 2	98	202	197	102
Game 3	123	180	238	134

Table 2: A table too wide to fit within page.

	Digits 1-33	Digits 34-66	Digits 67-99	Ref.
pi	3.14159265358979323 846264338327950	28841971693993751 0582097494459230	78164062862089986 2803482534211706	piday.org
e	2.71828182845904523 536028747135266	24977572470936999 5957496696762772	40766303535475945 7138217852516642	nasa.gov

Table 3: A table with merged cells using the `attributes` plugin.

	Colors	
Size	Text Color	Background Color
big	blue	orange
small	black	white

Equations

A LaTeX equation:

$$\int_0^{\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2} \quad (1)$$

An equation too long to fit within page:

$$x = a + b + c + d + e + f + g + h + i + j + k + l + m + n + o + p + q + r + s + t + u + v + w + x + y + z + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 \quad (2)$$

Special

⚠ WARNING The following features are only supported and intended for `.html` and `.pdf` exports. Journals are not likely to support them, and they may not display correctly when converted to other formats such as `.docx`.

LINK STYLED AS A BUTTON

Adding arbitrary HTML attributes to an element using Pandoc's attribute syntax:

Manubot Manubot Manubot Manubot Manubot. Manubot Manubot Manubot Manubot.
Manubot Manubot Manubot. Manubot Manubot. Manubot.

Adding arbitrary HTML attributes to an element with the Manubot `attributes` plugin (more flexible than Pandoc's method in terms of which elements you can add attributes to):

Manubot Manubot Manubot Manubot Manubot. Manubot Manubot Manubot Manubot.
Manubot Manubot Manubot. Manubot Manubot. Manubot.

Available background colors for text, images, code, banners, etc:

white lightgrey grey darkgrey black lightred lightyellow lightgreen
lightblue lightpurple red orange yellow green blue purple

Using the [Font Awesome](#) icon set:

✓ ? ★ 🔍 ✖ ...

 **Light Grey Banner**useful for *general information* - manubot.org **Blue Banner**useful for *important information* - manubot.org **Light Red Banner**useful for *warnings* - manubot.org

Introduction

Science journalism is an indispensable part of scientific communication and provides an accessible way for everyone from researchers to the public to learn about new scientific findings and to consider their implications. However, it is important to identify the ways in which its coverage may skew towards particular demographics. Coverage of science shapes who is considered a scientist and field expert by both peers and the public. This indication of legitimacy can either help recognize people who are typically overlooked due to systemic biases or intensify biases. Journalistic biases in general-interest, online and printed news have been observed by journalists themselves [9,10#44nchdhay,11,12/], as well as by independent researchers [13,14,15,16,17,18/]. Researchers found a gap between male and female subjects or sources, with independent studies finding that between 17-40% of total subjects were female across multiple general-interest printed news outlets between 1985 and 2015 [13,14,18/]. One study found 27-35% of total subjects in international science and health related news were female between 1995 and 2015, and 46% in print, radio, and television in the United States in 2015 [18/]. While gender disparities in news coverage have been extensively researched, research into disparities with respect to name origins is currently lacking in the literature.

It should be noted that scientific news coverage is confounded by the existing differences in gender and racial demographics within the scientific field [19,20]. However, we are interested in quantifying disparities with respect to observed demographic differences in the scientific field, using academic authorship as an estimate for the existing demographics. This is similar to other studies that have quantified gender or racial disparities in science as observed in citation [21,22] and funding rates [23,24,25,26,27].

In researching a story, a journalist will typically interview multiple sources for their opinion, potentially asking for additional sources, thus allowing individual unconscious biases at any point along the interview chain to skew scientific coverage broadly. In addition, the repeated selection of a small set of field experts or the approach a journalist takes in establishing a new source may intensify existing biases [28,11,12/]. While disparities in representation may go unnoticed in a single article, analyzing a large corpus of articles can identify and quantify these disparities and help guide institutional and individual self-reflection. In the same vein as previous media studies [13,14,15,16,17,18/], we sought to quantify gender and regional differences of journalism beyond the existing demographic differences in the scientific field. Our study focused solely on science journalism, specifically content published by *Nature*. Since *Nature* also publishes primary research articles, we used these data to determine the demographics of the expected set of possible sources. For clarity, throughout this manuscript we will refer to journalistic articles as *news* and academic, primary research articles as *papers*. Furthermore, when we refer to "authors" we mean authors of academic papers, not journalists. In our analysis, we identified quoted and cited people by analyzing the content and citations within all news articles from 2005 to 2020, and compared this demographic to the academic

publishing demographic by analyzing first and last authorship statistics across all *Nature* papers during the same time period.

Through our analysis of 22,001 news articles, we were able to identify >8,800 quotes and >4,000 citations with sufficient speaker or author information. We also identified first and last authors of >10,000 *Nature* papers. We then identified possible gender or regional differences using the extracted names. The extracted names were used to generate three data-types: quoted, mentioned, and cited people. We used computational methods to predict gender and identified a trend towards quotes from people predicted male in news articles when compared to both the general population and predicted male authorship in papers. Within the period that we examined, the proportion of predicted male attributed quotes in news articles went from initially higher to currently lower than the proportion of male first and last authors in *Nature* papers. Furthermore, we found that the quote difference was dependent on article type; the “Career Feature” column achieved gender parity in quoted speakers.

We also used computational methods to predict name origins of quoted, mentioned, and cited people. Through our analysis, we found a significant over-representation of names with predicted Celtic/English origin and under-representation of names with a predicted East Asian origin in both quotes and mentions. To our knowledge, our work is the first to identify a substantial under-representation of names with a predicted East Asian origin in scientific journalism.

While we focused on news from *Nature*, our software can be repurposed to analyze other text. We hope that publishers will welcome systems to identify disparities and use them to improve representation in journalism. Furthermore, our approach is limited by the features we were able to extract, which only reflects a portion of the journalistic process. Journalists could additionally track all sources they contact to self-audit. However, auditing is only part of the solution; journalists and source recommenders must also change their source gathering patterns. To help change these patterns, there exist guides [28/], databases [30], and affinity groups [28/] that can help us all expand our vision of who can be a field expert.

Methods

Data Acquisition and Processing

Text Scraping

We scraped all text and metadata from *Nature* using the web-crawling framework Scrapy [31/] (version 2.4.1). We created four independent scrapy web spiders to process the news text, news citations, journalist names, and paper metadata. News articles were defined as all articles from 2005 to 2020 that were designated as “News”, “News Feature”, “Career Feature”, “Technology Feature”, and “Toolbox”. Using the spider “target_year_crawl.py”, we scraped the title and main text from all news articles. We character normalized the main text by mapping visually identical Unicode codepoints to a single Unicode codepoint and stripping many invalid Unicode characters. Using an additional spider defined in “doi_crawl.py”, we scraped all citations within news articles. For simplicity, we only considered citations with a DOI included in either text or a hyperlink in this spider. Other possible forms of citations, e.g., titles, were not included. The DOIs were then queried using the *Springer Nature* API. The spider “article_author_crawl.py” scraped all articles designated “Article” or “Letters” from 2005 to 2020. We only scraped author names, author positions, and associated affiliations from research articles, which we refer to as *papers*. It should be noted that “News” article designations changed over time. Additionally, scraping for journalist names was performed months after the initial scraping of the text, and some aspects of the *Nature* website changed. Website changes caused us to

lose unique file mappings between the scraped journalist name and other article metadata for 137 articles. Less than thirty articles per year were impacted.

coreNLP

After the news articles were scraped and processed, the text was processed using the coreNLP pipeline [32] (version 4.2.0). The main purpose for using coreNLP was to identify named entities related to mentioned and quoted speakers. The full set of annotators were: tokenize, ssplit, pos, lemma,ner, parse, coref, quote. We used the “statistical” algorithm to perform coreference resolution. All results were output to json format for further downstream processing.

Springer Nature API

Springer Nature was chosen over other publishers for multiple reasons: 1) it is a large publisher, second only to Elsevier; 2) it covers multiple subjects, in contrast to PubMed; 3) its API has a large daily query limit (5000/day); and 4) it provided more author affiliation information than found in Elsevier. We generated a comparative background set for supplemental analysis with the *Springer Nature* API by obtaining author information for papers cited in news articles. We selected a random set of papers to generate the *Springer Nature* background set. These papers were the first 200 English language “Journal” papers returned by the *Springer Nature* API for each month, resulting in 2400 papers per year for 2005 through 2020. To obtain the author information for the cited papers, we queried the *Springer Nature* API using the scraped DOI. For both API query types, the author names, positions, and affiliations for each publication were stored and are available in “all_author_country.tsv” and “all_author_fullname.tsv”.

Name Formatting

Name Formatting for Gender Prediction in Quotes or Mentions

We first pre-filter articles that have more than 25 quotes, which is 2.69% (433/16,080) of total articles. This was done to ensure no single article is over-represented and to avoid spuriously identified quotes due to unusual article formatting. To identify the gender of a quoted or mentioned person, we first attempt to identify the person’s full name. Even though genderizeR only uses the first name to make the gender prediction, identifying the full name gives us greater confidence that we are using the first name. To identify the full name, we take the predicted speaker by coreNLP and match it to the longest matching name within the same article. We match names by finding the longest mentioned name in the article with minimal edit (Levenshtein) distance. The name with the smallest edit distance, where character deletions have zero cost, is defined as the matching name. Character deletion was assigned a zero cost because we would like exact substring matches. For example, the calculated cost, including a cost for character deletion, between John and John Steinberg is 10; without character deletion, it is 0. Compared with the distance between John and Jane Doe, with character deletion cost, it is 7; without it is 2. If we are still unable to find a full name, or if coreNLP cannot identify a speaker at all, we also determine whether or not coreNLP linked a gendered pronoun to the quote. If so, we predict that the gender of the speaker is the gender of the pronoun. We ignore all quotes with no name or partial names and no associated pronouns. A summary of processed gender predictions of quotes at each point of processing is provided in Table 4.

Name Formatting for Gender Prediction of Authors

Because we separate first and last authors, we only considered papers with more than one author. As for quotes, we needed to extract the first name of the authors. We cast names to lowercase and processed them using the R package humaniformat [33]. humaniformat is a rule based program that

uses character markers to identify if names are reversed (Lastname, Firstname), find middle names and titles. This processing was not required for quote prediction because names written in news articles did not appear to be reversed or abbreviated. Since many last or first authorships may be non-names, we additionally filtered out any identified names if they partially or fully match any of the following terms: “consortium”, “group”, “initiative”, “team”, “collab”, “committee”, “center”, “program”, “author”, or “institute”. Furthermore, since many papers only contain first name initials (for example, “N. Davidson”), we remove any names less than four letters (length includes punctuation) and containing a “.” or “-”, then strip out all periods from the first name. This ensures that hyphenated names are not changed, e.g. Julia-Louise remains unchanged, but removes hyphenated initials, e.g. J-L. A summary of processed author gender predictions at each point of processing is provided in Tables 5 - 7.

Name Formatting for Name Origin Prediction

In contrast to the gender prediction, we require the entire name in all steps of name origin prediction. For names identified in the *Nature* news articles, we use the same process as described for the gender prediction; we again try to identify the full name. For author names, we process the names as previously described for the gender prediction of authors. For all names, we only consider them in our analyses if they consist of two distinct parts separated by a space and excluding titles (e.g. Mrs., Prof., Dr., etc.). All names that were filtered out in the analysis of quotes and mentions are provided on our github in the file “data/author_data/all_mentioned_fullname_excluded.tsv” and “data/author_data/all_speaker_fullname_excluded.tsv”. A summary of processed name origin predictions of quotes and citations at each point of processing is provided in Tables 4 - 7.

Gender Analysis

The quote extraction and attribution annotator from the coreNLP pipeline was employed to identify quotes and their associated speakers in the article text. In some cases, coreNLP could not identify an associated speaker’s name but instead assigned a gendered pronoun. In these instances, we used the gender of the pronoun for the analysis. The R package genderizeR [34], a wrapper for the genderize.io API [35/], predicted the gender of authors and speakers. We predicted a name as male using the first name with a minimum cutoff of 50%. To reduce the number of queries made to genderize.io, a previously cached gender prediction from [36] was also used and can be found in the file “genderize.tsv”. All first name predictions from this analysis are in the file “genderize_update.tsv”. To estimate the gender gap for the quote gender analyses, we used the proportion of total quotes, not quoted speakers. We used the proportion of quotes to measure speaker participation instead of only the diversity of speakers. The specific formulas for a single year are shown in equations 3 and 4. We did not consider any names where no prediction could be made or quotes where neither speaker nor gendered pronoun was associated.

$$\text{Prop. Male Quotes} = \frac{|\text{Male Speaker Quotes}|}{|\text{Male or Female Speaker Quotes}|} \quad (3)$$

$$\text{Prop. Male First Authors} = \frac{|\text{Male First Authors}|}{|\text{Male or Female First Authors}|} \quad (4)$$

Name Origin Analysis

We used the same quoted speakers as described in the previous section for the name origin analysis. In addition, we also consider all authors cited in a *Nature* news article. In contrast to the gender prediction, we need to use the full name to predict name origin. We submitted all extracted full names to Wiki-2019LSTM [36] to predict one of ten possible name origins: African, Celtic/English, East Asian, European, Greek, Hispanic, Hebrew, Arabic/Turkish/Persian, Nordic, and South Asian. While a full description of Wiki-2019LSTM is outside the scope of this paper, we describe it here briefly. Wiki-2019LSTM is trained on name and nationality pairs, using 3-mers of the characters in a name to predict a nationality. To ensure robust predictions, nationalities were grouped together as described in NamePrism [37]. NamePrism chose to exclude the United States, Australia, and Canada from their country groupings and were therefore excluded during training of Wiki-2019LSTM. This choice was justified by NamePrism in stating that these countries had a high level of immigration. The treemap of country groupings defined in the NamePrism manuscript are found in figure 5 of the publication [37].

After running the pre-trained Wiki-2019LSTM model, we select the highest probability origin for each name as the resultant assignment. Similar to the gender analyses, quote proportions were again directly compared against publication rates. For citations, quotes, and mentions, we calculated the proportion for a given year for each name origin. This is shown in 5 to, for example, calculate the citation rate for last authors with a Greek name origin for a single year.

$$\text{Prop. Greek Last Author Cited} = \frac{|\text{Cited Last Authors w/Greek Name}|}{|\text{Cited Last Authors w/any Name}|} \quad (5)$$

$$\text{Prop. Greek Quotes} = \frac{|\text{Quotes w/Greek Named Speaker}|}{|\text{Quotes w/any Named Speaker}|} \quad (6)$$

$$\text{Prop. Greek Names Mentioned} = \frac{|\text{Unique Greek Names Mentioned}|}{|\text{Unique Names w/any Origin Mentioned}|} \quad (7)$$

Identifying Quotes or Mentions with US Affiliation

We assigned affiliations to quoted or mentioned people when their name was also a cited first or last author in the same news article. All country affiliations within the cited article were assigned to the quoted or mentioned person. For example, if a researcher was affiliated with an Austrian university, but the cited paper has authors from both Austria, France, and the United States, the researcher will be given three affiliations.

Bootstrap Estimations

For all analyses related to equations 3 - 7, we independently selected 5000 bootstrap samples for each year. We sampled with replacement of size equal to the cardinality of the complete set of interest. Bootstrap estimates for equations 3 - 7 were performed by sampling the denominator set. The mean, 5th, 95th quantiles across the estimates are reported as the estimated mean, lower, and upper bounds. For the divergent word analysis, due to computational constraints, we only took 1000 bootstrap samples. The bootstrap estimates were taken by subsampling the news articles with replacement, each time recalculating the country-normalized token frequencies within each country set (C and M). After the normalized frequencies within each country set were calculated, we calculated the ratio between country sets for each subsample with a pseudocount of 1 in the numerator and

denominator, $(C+1)/(M+1)$. Again, the mean, 5th, 95th quantiles across the estimates are reported as the estimated mean, lower, and upper bounds.

Results

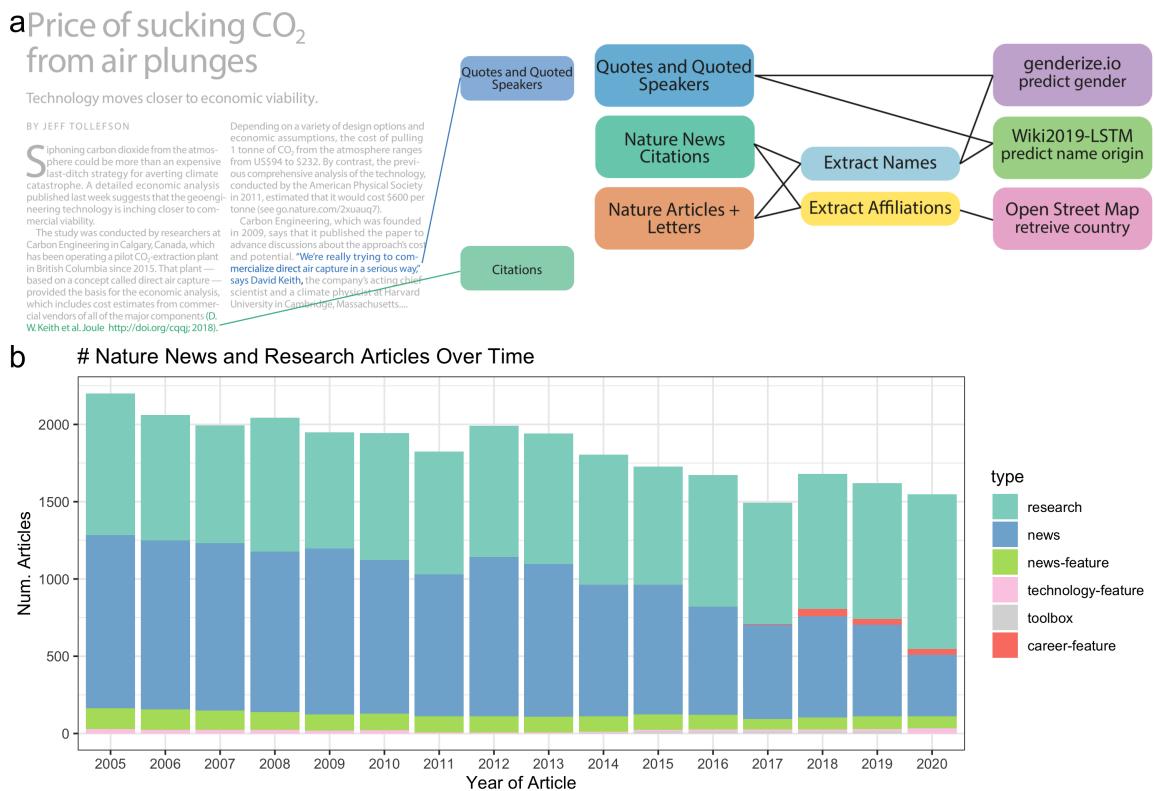


Figure 5: Data and Processing Pipeline Overview Panel A, left, depicts an example news article and the type of data extracted from the text. Green and blue highlighted text depicts all quotes, and associated speakers identified by the coreNLP pipeline. A custom script described in section [Methods](#) identifies all citations. Panel A, right, charts the analyses done on the extracted names and locations from news articles and papers published by *Nature*. Panel B shows the types and amounts of articles that we have used for analyses.

Creation of an Annotated News Dataset

We have analyzed the text of 22,001 news-related articles hosted on “www.nature.com” that span 15 years from 2005 to 2020. Our primary focus is on 16,080 articles written by journalists which include the following five article types: “Career Feature”, “News”, “News Feature”, “Technology Feature”, and “Toolbox”. “Career Feature” generally focuses on the career-related aspects of being a scientist. “News” and “News Feature” focuses on current events related to science as well as new scientific findings. It should be noted that the types of articles contained in “News” changed over time which may induce content shifts in a subset of the articles within our corpus. “Technology Feature” also covers current events and scientific findings, but additionally focuses on how science intersects with technology, such as apps, methodologies, tools, and practices. Lastly, “Toolbox” is similar to “Technology Feature”, but is more centered on technology, especially the tools used to perform science. We also include one analysis of the scientist-written news articles, “Career Column” and “News and Views”, as an additional set of 5,921 articles. “Career Column” is similar to “Career Feature”, except it is not written by journalists, but individuals in the scientific field. “News and Views” is similar to a review article, where a field expert writes an article relating to a recently written article within *Nature*.

The top three observed article frequencies are “Research” (including “Letters” and “Articles”), “News”, and “News Feature”. Since *Nature* merged “Letters” and “Research” papers in 2019, we combined

them in our analysis. We observed substantial variability in the number of *Nature* news articles by type between 2005 and 2020 (Figure 5b). The changing classification of article types may explain temporal changes in news articles. Over time, the frequency of “News” articles decreased; however, more specific news-related article types increased, including the introduction of the new categories “Career Feature”, “Toolbox”, and “Career Column”.

Terms used in Analysis

The text and citations were then uniformly processed as depicted in Figure 5a to identify: 1) quotes and quoted speakers (blue box) and 3) cited authors (green box). The extracted names from the text were used to generate three data types for downstream processing: quoted, mentioned, and cited people. A summary of frequencies for each data type at each point of processing is provided in Tables 4 - 7. We scraped the text using the web-crawling framework Scrapy [31/], processed, and ran it through the coreNLP pipeline ([Methods](#)). To identify quotes and speakers, we used the coreNLP quote extraction and attribution annotator. We performed multiple name formatting processes ([Methods](#)) to identify the speaker’s full name for gender and name origin prediction. All names where we could identify two name parts, assumed to be a first and last names exclusive of titles, were used for gender prediction and checked against the genderize.io database. Since names used in the name origin analysis were computationally analyzed and not checked for existence in an existing database, we used additional filters ([Methods](#)).

All names excluded from the name origin analysis of quotes and mentions are provided on our github in the file “data/author_data/all_mentioned_fullname_excluded.tsv” and “data/author_data/all_speaker_fullname_excluded.tsv”. We found that most names were excluded because two name parts, assumed to be a first and last names exclusive of titles, were not found. We scraped the citations using an independent scraper to the text scraper, but still utilizing the Scrapy framework. All identified DOI’s were queried using the *Springer Nature* API to attain all authors’ names, positions, and affiliations. Country of affiliation was determined using OpenStreetMap [38], a free and widely used geographic database.

Comparator Datasets

Next, we determined if the quoted speakers and cited authors in news articles have a similar demographic makeup as the scientists who publish their primary research in *Nature*. To make this determination, we used all authors’ names, positions, and affiliations of papers published by *Nature* over the same time period (Figure 5a, dark orange box). The author metadata of *Nature* papers from 2005 to 2020 totaled 13,414. To more broadly represent overall science authorship, we also separately analyzed 38,400 randomly selected *Springer Nature*-published papers from English-language journals over the same time. It should be noted that extracted quotes may come from multiple types of people, such as academic scientists, clinicians, the broader scientific community, politicians, and more. However, through anecdotal observation we believe that most sources come from either academic scientists or those actively involved in science. The extracted author affiliations from both data sources were mapped to a country using OpenStreetMap. Similarly, author names were uniformly processed and then used to predict both gender and name origin.

Quoted Speakers and Primary Research Authors in *Nature* are More Often Male

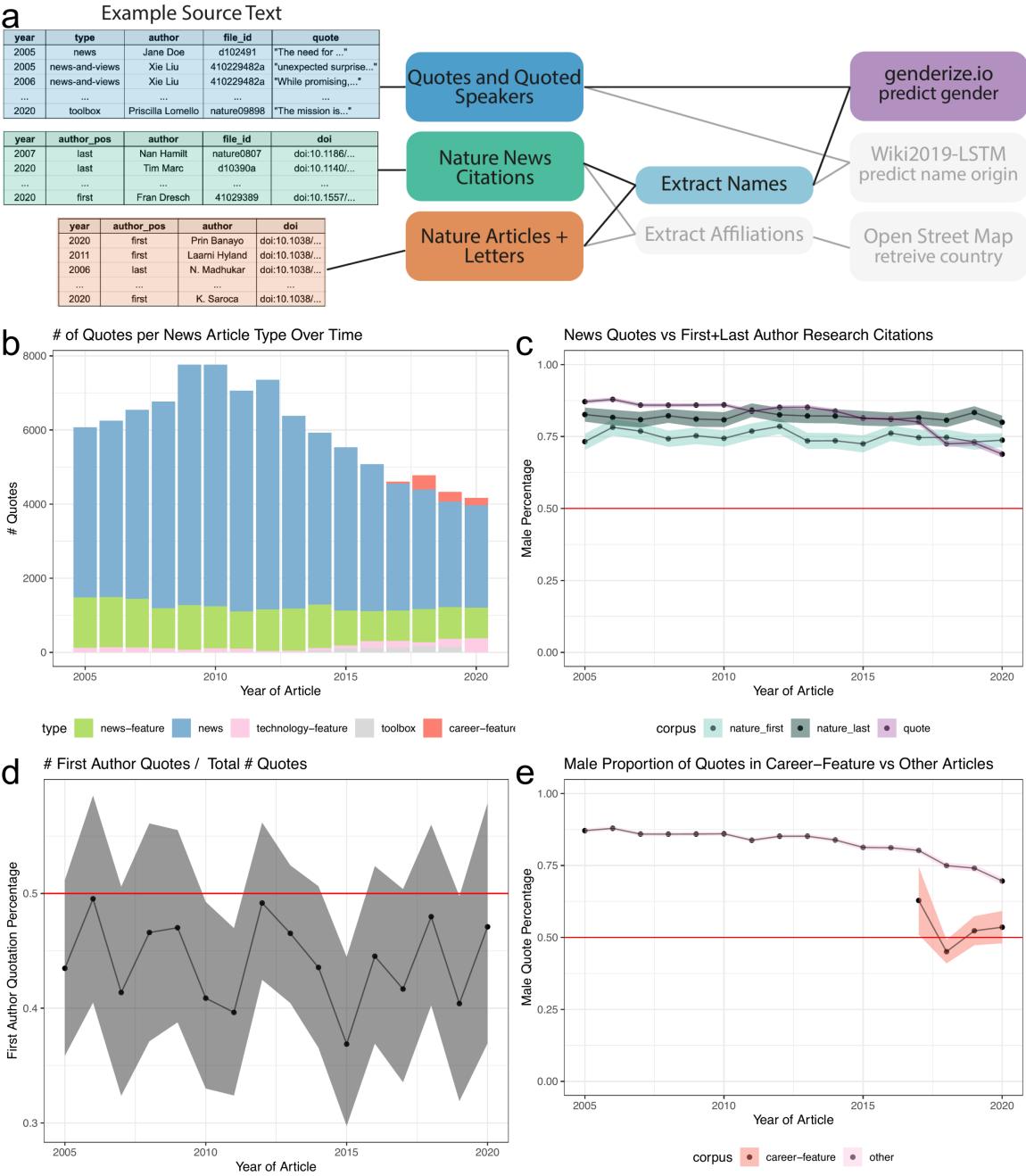


Figure 6: Predicted male speakers are overrepresented in quotes, but this depends on the article type. Panel A, left, depicts an example of the names extracted from quoted speakers in news articles and authors in papers. Panel A, right, highlighted the data types and processes used to analyze the predicted gender of extracted names. Panel B shows an overview of the number of quotes extracted for each article type. Panel C depicts three trend lines: Purple: Proportion of quotes for an estimated male speaker; Light Blue: Proportion of first author papers from an estimated male author; Dark Blue: Proportion of predicted male last authors. We observe that the proportion of estimated male quotes is steadily decreasing, most notably from 2017 onward. This decreasing trend is not due to a change in quotes from the first or last authors, as observed in Panel D. Panel D shows a consistent but slight shift towards quoting the last author of a cited article than the first author. Instead, the observed downward trend of male quotes coincides with additional article types introduced in 2017. Panel E depicts the frequency of quote by article type highlighting an increase in quotes from "Career Feature" articles. Panel E depicts that the quotes obtained in this article type have reached parity. The colored bands represent a 5th and 95th bootstrap quantiles in all plots, and the point is the mean calculated from 5,000 bootstrap samples.

To quantify and compare the gender demographic of quoted people and authors, we analyzed their names. While we could have analyzed the proportion of unique male speakers, we were interested in measuring the overall participation rates by gender and analyzed the proportion of total quotes, e.g. a single speaker may have more than one quote in an article. Furthermore, we assume that a majority of quoted speakers are typically involved in scientific research and therefore primary research authors is a comparable demographic. Figure 6 shows an overview of the process and example input

data for this analysis: 1) quotes and quoted speakers (blue box), 2) first and last authors' names of papers published by *Nature* (dark orange box). These analyses relied upon accurate gender prediction of both authors and speakers. To predict the gender of the speaker or author, we used the package genderizeR [34], an R package wrapper to access the genderize.io API [35/] to get binary gender predictions for each identified first name. We unfortunately cannot identify non-binary gender expression with the tools we used. Performance of binary prediction was evaluated on a benchmark data set of thirty randomly selected news articles, ten from each of the following years: 2005, 2010, 2015 (Figure [Supplemental 1a](#)). In addition, genderize.io has been found by independent researchers to have an error rate comparable to other published gender prediction methods, with a error-rate on predicted names below 6% [39,40]. However, it should be noted that the error rate varies by name origin with the largest decrease in performance on names with an Asian origin [39,40] .

We first examined the number of quotes identified within each type of science-news article (Figure [6b](#)), totaling 105,457 quotes with 96,390 of them containing a gender prediction for the speaker. Quote frequencies vary by article type. We compared the number of quotes from predicted male people to the number of predicted male first and last authors published in *Nature*. The total number of authors with a gender prediction were 10,601 first authors and 10,572 last authors across a total of 11,161 publications. As denoted by the red line, we found that the predicted genders of authors and source-quotes were far from gender parity (Figure [6c](#)). We found this result consistent for articles written by either a predicted male or female journalist (Figure [Supplemental 2a,b](#)). Additionally, we observed a difference in the predicted genders between first and last authors, with the last authors more frequently predicted to be male.

Comparison with *Springer Nature*

To extend our analysis to primary research authors more broadly, we also examined a random selection of authors from English language journals published by *Springer Nature* (Figure [Supplemental 3a](#)). The predicted gender gap between first and last authors was larger in our selection of *Springer Nature* papers; however, both first and last authors were predicted to be closer to parity than for *Nature* authors. Overall, predicted male people were more frequently quoted than predicted female people in *Nature* news articles and first and last authors in *Nature* and *Springer Nature* papers over the same time period.

Career Feature Articles reach Gender Parity

The gender proportions of authorship were relatively stable over time for both *Nature* and *Springer Nature* papers. In contrast, we found that the rate of quotes predicted to be from male people noticeably decreased over time. In 2005, the fraction of quotes predicted to be from male people was 87.09% (5,291/6,075) whereas in 2020 it was 68.86% (2,870/4,168). We identified that a large decrease occurred in *Nature* between 2017 and 2018. We explored the possible reasons for this decrease. First, we looked at the authorship position of speakers who were quoted about their published paper (Figure [6d](#)). We identified 6,545 quotes with an associated citation (2,871 first author and 3,674 last author quotes). We found that quotes trend slightly towards last authors from 2005 to 2020, but because the fraction of predicted male last authors remained stable over time both for *Nature* and the selection of *Springer Nature* papers, which likely does not explain the downward trend. We then analyzed the breakdown of gender predicted quotes by article type. Interestingly, one article type, "Career Feature", achieved gender parity in its quotes (Figure [6e](#) and Figure [Supplemental 3b](#)). In this article type, we identified a total of 898 quotes (449 predicted female and 449 predicted male quotes), which substantially pulled the overall quote gender ratio closer to parity from 2018 onward.

Predicted Celtic English Name Origins are over-enriched in quoted and mentioned people, while predicted East Asian name origins are under-

enriched

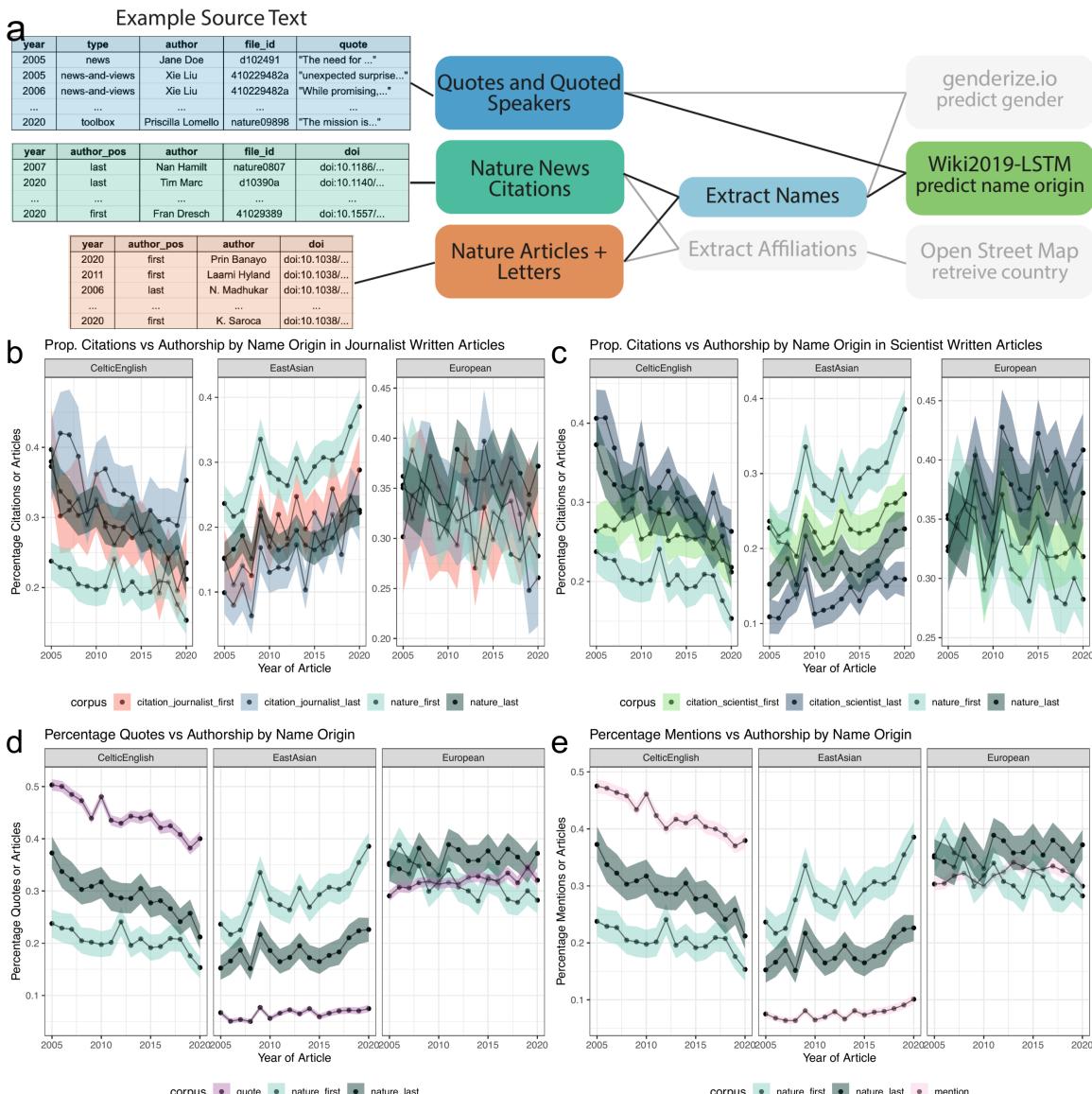


Figure 7: Analysis of Quotes and Citations found Over-representation of Celtic/English and under-representation of East Asian predicted name origins. Panel A, left, depicts an example of the names extracted from quoted speakers and citations found within news articles and authors in papers. Panel A, right, highlights the data types and processes used to analyze the predicted origin of extracted names. Panels B and C depict a comparison between the predicted name origins of last authors in *Nature* and cited papers in the news. Panel B and C differ in the news article types. Panel B calculates the predicted name origin proportion using only journalist-written articles, whereas Panel C only uses scientist-written articles. The distinction between scientist- and journalist-written articles are defined by the article appearing in either the "Career Column" or "News and Views" sections, or another section, respectively. Similarly, Panels D and E depict two possible trend lines, comparing predicted name origins of either quoted or mentioned people against name origins of last authors of *Nature* research papers. .

To identify possible disparities with respect to name origin, we again used the extracted names of quoted speakers from *Nature* news articles and last authors of published papers in *Nature*. In addition, we also identified the last authors of all papers cited by a *Nature* news article. All processed names were then input into Wiki2019-LSTM and assigned one of ten possible name origins ([Methods](#)). Figure 7a shows an overview of the process and example input data for this analysis: 1) quotes and quoted speakers (blue box), 2) names of cited first and last authors in news articles (green) 3) first and last authors' names of papers published by *Nature* (dark orange box). We divided our analysis into three parts: firstly, quantifying the proportions of predicted name origins of first and last authors cited in *Nature* news articles. Secondly, calculating the proportion of quotes from speakers with a predicted name origin. Thirdly, calculating the proportion of unique names mentioned within an article with a

predicted name origin. As a comparator set, we again used the first and last author names in *Nature* papers for all three analyses. Additionally, in our supplemental analyses, we compared against the first and last authorship in a random selection of *Springer Nature* papers. We found that the number of quotes and unique names mentioned dramatically outnumbered the number of cited authors in *Nature* news articles, as well as first and last authors within *Nature* papers (Figure [Supplemental 4a](#)). Still, since we have more than one hundred observations per time point for each data type, we believe this is sufficient for our analysis. Minimum and median per data type over all years: *Nature* papers, (568, 684); *Springer Nature* papers, (1332, 1710); *Nature* quotes, (3788, 5696); *Nature* mentions, (3225, 4752); citations in journalist-written *Nature* article, (142, 268) citations in a scientist-written *Nature* article, (512, 664).

News citation rates across name origin predictions nearly match *Nature* authorship

In comparing the citation rate of first and last author name origins in news articles, we decided to additionally analyze scientist-written articles. Though fewer in number, scientist-written news articles have many citations, making the set sufficient for analysis and providing an opportunity to measure differences in citation patterns between journalists and scientists. In both journalist- and scientist-written articles, we found that most cited name origins were predicted Celtic/English or European, both with a bootstrapped estimated citation rate between 19.2-42.8% (Figure [Supplemental 4b,c](#)). East Asian predicted name origins are the third highest proportion of cited names, with a bootstrapped estimated citation rate between 6.4-28.8%. All other predicted name origins individually account for less than 7.9% of total cited authors.

We analyzed how these distributions compare to the composition of the first and last authors in *Nature* (Figure [Supplemental 5](#)), by examining the top three most frequent predicted name origins (Figure [7b,c](#)). When considering only first authors, we found a slight over-enrichment for predicted Celtic/English name origins and a small under-enrichment for predicted East Asian name origins in scientist-written and journalist-written news articles when compared to the composition of first authors in *Nature* (Figure [7b, c](#)). When considering last authors, this pattern no longer exists. Furthermore, we found no substantial difference for European or other predicted name origins when comparing against first and last authorship within *Nature* (Figure [Supplemental 6a](#)). We also observed the predicted Celtic/English over-enrichment and East Asian under-representation when considering our subset of *Springer Nature* papers (Figure [Supplemental 6b](#)) for both journalist- and scientist-written news articles. In contrast to *Nature*, in the *Springer Nature* set, we see a difference in predicted European name origins, with a growing over-enrichment. Additionally, we see a difference in predicted Arabic/Turkish/Persian name origins frequencies between cited authors and *Springer Nature* authors, however the absolute difference is lower than observed for Celtic/English and East Asian predicted name origins.

News quotation rates are over-enriched for predicted Celtic-English and under-enriched for East Asian name origins.

We then sought to determine whether or not the quoted speaker demographic replicated the cited authors' over- and under-enrichment patterns. We found a much stronger Celtic/English over-enrichment in comparison to citation patterns, with quotes from those with Celtic/English name origins at a much higher frequency than quotes from those with European name origins (Figure [Supplemental 4d](#)). We also found a much stronger depletion of quotes from people with predicted East Asian name origins (Figure [Supplemental 4b](#)), with never more than 7.7% of quotes (Figure [7d](#)). This reveals a large disparity when considering that people with a predicted East Asian name origin constitute between 14.3-33.6% of last authors cited in either journalist- or scientist-written news articles (Figure [7b,c](#)). When we compare against first and last authorship in *Nature* across all predicted

name origins, we find that for all other name origins except for East Asian and Celtic/English, the quote rates closely matches the predicted name origin rate of first and last authors in *Nature* (Figure [Supplemental 6c](#), dark grey and light blue lines compare to the purple lines).

To further understand the source of Celtic/English over-enrichment and East Asian under-enrichment, we selected a subset of quotes from people whose works were also cited in the news article. The purpose of this additional comparison of quoted speakers versus quoted *and* cited speakers was to reveal source gathering patterns beyond cited works. We found that the proportion of predicted East Asian name origins was closer to the expected rate after considering only quoted speakers with citations, more closely matching the analysis on citations alone (Figure [Supplemental 7a,b](#)). This indicates that expert opinions gathered beyond manuscript authors is responsible for a large proportion of the observed name disparities.

Next, we designed an experiment to test if predicted journalist name origin had any effect on quote disparities. We found that journalists with a predicted East Asian name origin had a higher rate of East Asian quoted speakers (24.3%) in comparison to journalists with Celtic/English (3.8%) or European (8.6%) predicted name origins (Table [8](#)). To examine if this was again driven by source gathering beyond manuscript authors we again subsetted the quotes by adding two constraints: 1) the quotes must be from a cited first or last author in the same news article (Table [9](#)) and 2) that the cited article must have a US affiliation (Table [10](#)). We found that differences between journalists with different predicted name origins was nearly eliminated when restricting to quoted and cited speakers, and with the additional restriction of US affiliated citations, as evidenced in the predicted East Asian column of Table [10](#). The differences between Table [8](#) and Tables [9](#) and [10](#), indicate that the predicted name origin of a journalist has some association with sources gathered outside of directly cited works.

When comparing *Nature* articles against the *Springer Nature* set of first or last authors, we again find the same patterns in quoted speakers with East Asian, Celtic/English, and Arabic/Turkish/Persian predicted name origins when comparing against the as we did in the previous citation analysis (Figure [Supplemental 6d](#), green and purple lines). In addition, we find an under-enrichment of predicted Hispanic, South Asian, and Hebrew name origins when comparing against the predicted name origin rate of first and last authors in our *Springer Nature* set.

News mention rates are over-enriched for predicted Celtic-English and under-enriched for East Asian name origins.

Since many journalists use additional sources that are not directly quoted, we also analyzed likely paraphrased speakers, e.g. a case in which the person was a source and mentioned in the story but not directly quoted. To do this, we identified all unique names that appeared in an article, which we term *mentions*. We found the same pattern of over-enrichment for predicted Celtic/English name origins and under-enrichment for East Asian name origins when comparing against both *Nature* and *Springer Nature* first and last authorships (Figure [7e](#), Figure [Supplemental 4d,e](#), Figure [Supplemental 6e,f](#)). Similar to the quote analysis, we selected a subset of mentions from people that were also cited in the news article. We again found that the disparity was greatly reduced (Figure [Supplemental 7c,d](#)).

Discussion

Science journalism is the critical conduit between the academic and public spheres, and consequently shapes the public's view of science and scientists. However, as observed in other forms of recognition in science, biases may shift coverage away from the known demographics within science [\[41\]](#). Ideally, scientific journalism is representative of academic papers. Though it would be best for news coverage to promote equitable representation, at a minimum quotes and citations would ideally match the regional and gender demographics of scientific academia. To examine this last point, we analyzed

22,001 news articles published in *Nature*, to identify quoted, mentioned, and cited people. We then compared this to the authorship statistics from *Nature's* papers and a subset of *Springer Nature's* English language papers.

We first looked at possible gender differences in quotes and found in both news outlets, a large, but decreasing, gender gap when compared to the broader population in all but one article type. Additionally, this result was consistent in articles written by both predicted female and male journalists. We found that the decreasing trend in *Nature* articles was largely driven by the recent introduction of a single column, "Career Feature". This column has an equal number of quotes from both genders, showing that gender parity is possible in science journalism. This finding, coupled with the near equal number of article written by male and female predicted journalists, argues for more diversity in topical coverage. Including more content that is not primarily focused on recent publications, but all topics surrounding the practice of science, may help to rapidly achieve gender parity in journalistic recognition. However, we do recognize that different journalistic columns have different purposes or may represent different demographics and be inherently more difficult to reach parity.

To further our analysis of possible coverage disparities, we looked to differences in predicted name origins of quoted and cited authors across all the processed news articles. Our findings provide additional support for previous studies that identified under-citation [42] and under-recognition [41] of East Asian people. Interestingly, we found under-citation of people with predicted East Asian name origins to be much less pronounced than under-quotation. We do not believe that the under-quotation is driven by paraphrasing sources, which may occur more frequently with non-native English speakers. We also found that the disparity observed in quotes and mentions was almost eliminated when only considering people that were additionally cited within the same article. This suggests that the source of the disparity may lie in the search for additional expert opinions.

Either way, the clear disparity of predicted East Asian researcher quotes and mentions argues for including a broader set of voices when seeking opinions beyond the academic papers being covered in the article. One solution could be to have region-specific journalists. While we were not directly able to examine the regions journalists lived in, this potential strategy is supported by our analysis of journalists with a predicted East Asian name origin. When considering quotes from people with a predicted East Asian name origin, we found that journalists who themselves have a predicted East Asian name origin include a higher proportion of these quotes than journalists with European or Celtic/English predicted names. When considering only people who were both quoted and cited, the effect of the predicted name origins of journalists was substantially dampened. We are unable to identify if this is a geographic bias of the reporters in this analysis, since we do not know the location of the journalist at the time of writing the article. However, having reporters explicitly focused on specific regional sources to better cover international opinions in science can help ameliorate this disparity.

Through our comprehensive analysis, we were able to quantify how recognized persons in news journalism vary by name origin and gender, then compare it to scientific publishing background rates. While we found a significant gender disparity, the rate of female representation in scientific news is increasing and outpacing first and last authorships on scientific papers. Furthermore, we identified a significant depletion of quotes from scientists with a predicted East Asian name origin when compared to paper authorship, and a significant but smaller depletion of cited authors with a predicted East Asian name origin in news content.

Previous anecdotal studies from journalists have shown that awareness of their bias can help them to reduce it [10,11/12]. Once a bias is identified an individual can seek resources to help them find and retain diverse sources, such as utilizing international expert databases like gage [29/] and SheSource

[30/]. Additional tips for journalists to achieve and maintain a diverse source pool is described by Christina Selby in the Open Notebook [28/].

It should also be mentioned that we were only able to analyze the data provided through scraping "www.nature.com". This is a major limitation, because the only measures that we have of demographics of sources are people who have their name mentioned or research cited within the article. Journalists do not quote or mention all of the sources that they interviewed or cite all of the papers that they read when researching an article. For example, a person may not be mentioned or quoted in the article because of length limitations, because they do not want to be named, or if they provide information that is not directly quotable but that still shapes the content of the article. A more accurate reflection of journalists' sources would be a self-maintained record of people they interview. Our work examines disparities with respect to recognition within articles, which can be measured by mentions, quotes, or citations of people.

Furthermore, many journalists are limited by who responds to their requests for an interview or recommendations from prominent scientists. Scientists fielding reporter inquiries can also audit themselves to examine the extent to which there are disparities in the sets of experts they recommend. Journalists and the scientists they interview have a unique opportunity to shape the public and their peers' perspectives on who is a scientific expert. Their choice of coverage topics and interviewees could help to reduce disparities in the outputs of science-related journalism.

Data and Resource availability

This manuscript was written using Manubot [8] and is available on github: [manuscript repository link](#). All code and metadata is also available on github, [full analysis repository link](#), under a BSD 3-Clause License. The code to generate all main and supplemental figures are available as R markdown documents within our main analysis github, in the following subfolder: [notebooks](#). Due to copyright, we are unable to provide the scraped data used in this analysis. However, scraping code is available on our main analysis github, in the following subfolder: [scraper](#). To ensure reproducibility without violating copyright, we provide the word frequencies for each news article and the coreNLP output. Furthermore, we provide a docker image that can re-run the analysis pipeline using intermediate, pre-processed data and produce all the main and supplemental figures. To re-run the entire pipeline (including scraping), the docker image contains all necessary packages and code. The shell scripts to re-run the entire analysis are provided in the README file in the github repository.

Acknowledgements

We would like to thank Jeffrey Perkel for asking thoughtful questions that spurred this line of research, and providing feedback and insight into the news-gathering process during the course of this project.

Supplemental Figures

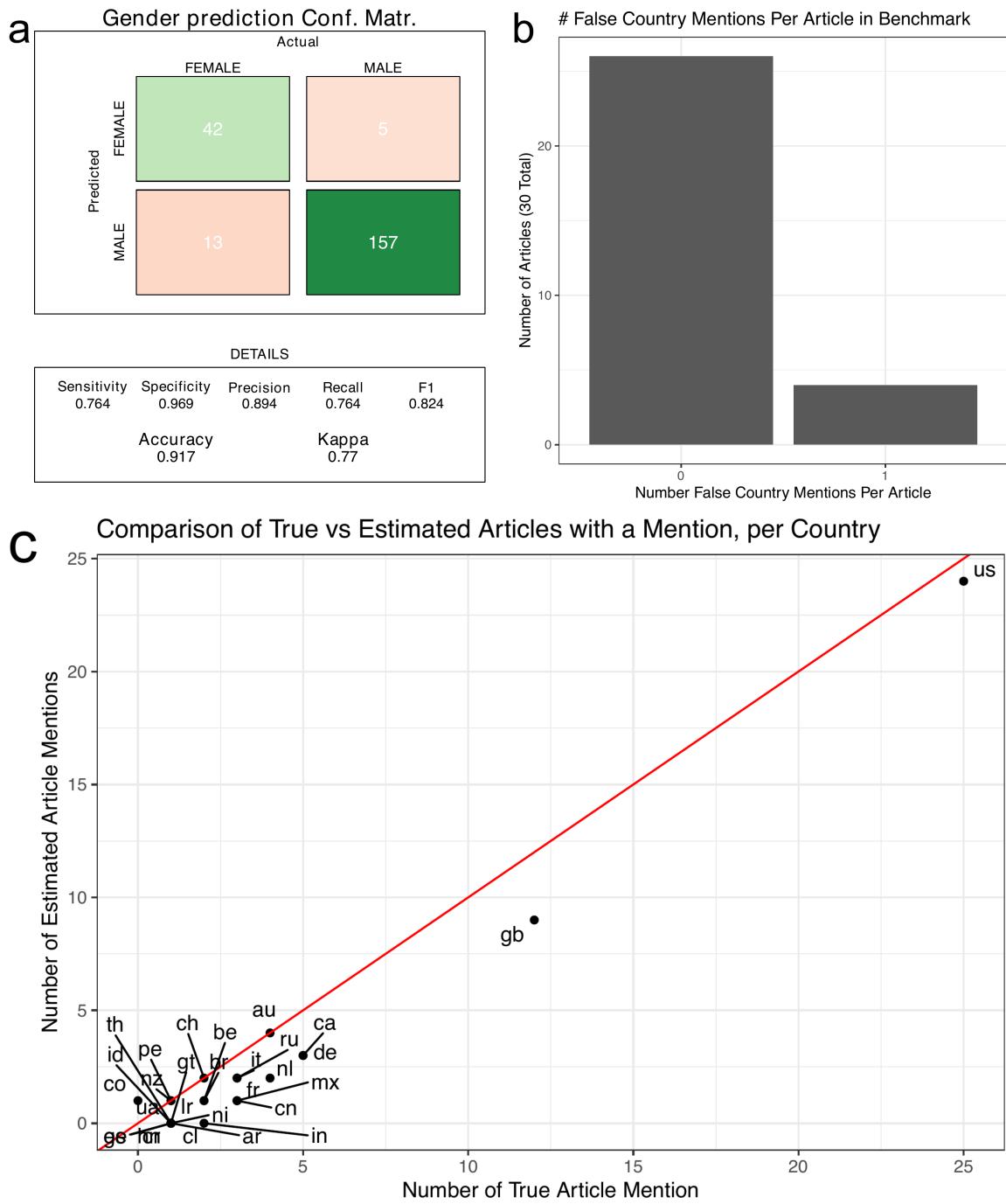


Figure Supplemental 1: Benchmark Data Panel A, depicts the performance of gender prediction for pipeline-identified quoted speakers. Panel B is a histogram of the number of articles that were falsely identified to mention a country by our processing pipeline. Panels C shows the estimated versus true frequency of country mentions within our benchmark dataset. The red line denotes the $x = y$ line.

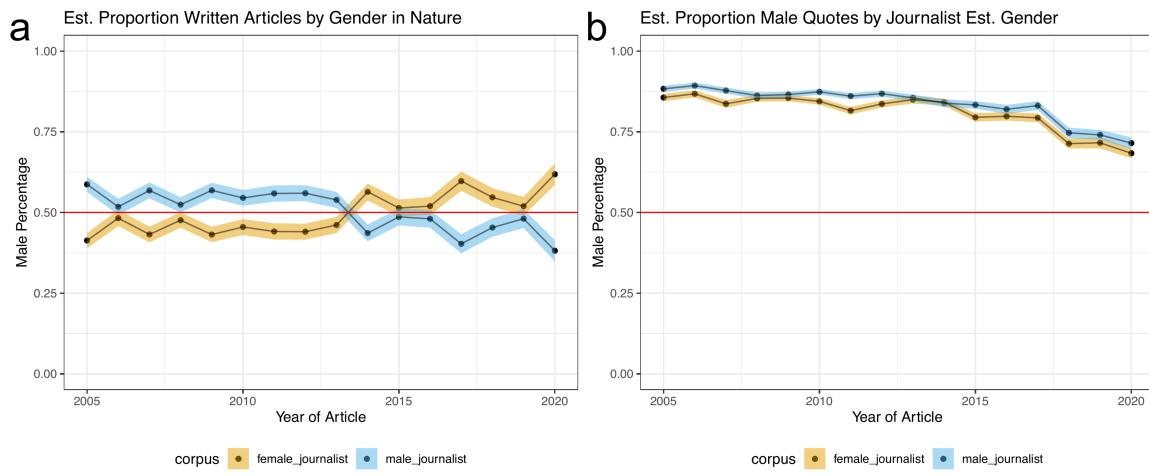


Figure Supplemental 2: Predicted male speakers are overrepresented in news quotes regardless of predicted journalist gender Panel A depicts two trend lines: Yellow: Proportion of *Nature* news articles written by a predicted female journalist; Blue: Proportion of *Nature* news articles written by a predicted male journalist. We observe almost no gender difference in the number of articles written by male and female journalists. Panel B depicts two trend lines: Yellow: Proportion of predicted male quotes in an article written by a predicted female journalist; Blue: Proportion of predicted male quotes in an article written by a predicted male journalist. In all plots, the colored bands represent the 5th and 95th bootstrap quantiles and the point is the mean calculated from 5,000 bootstrap samples.

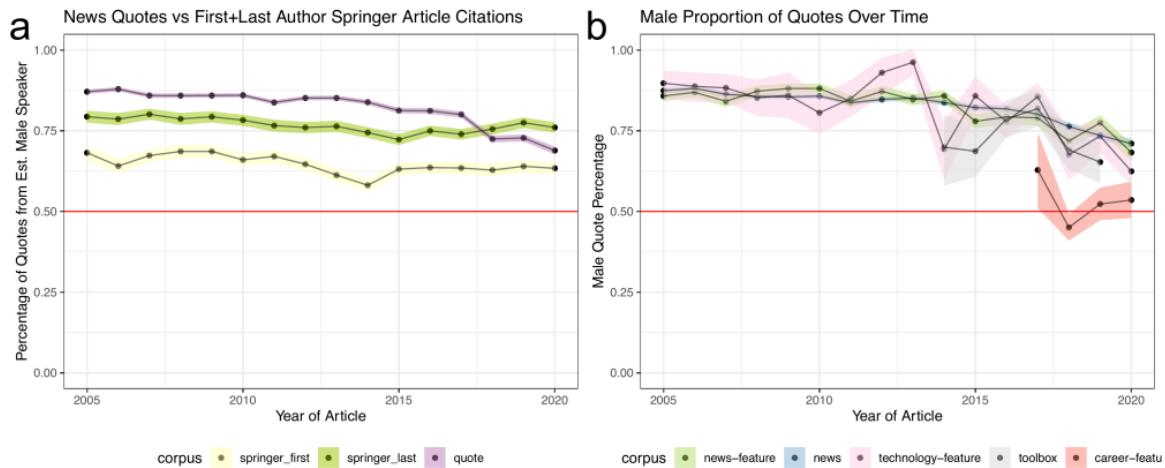


Figure Supplemental 3: Predicted male speakers are overrepresented in news quotes when compared against Springer Nature authorship Panel A depicts three trend lines: Purple: Proportion of *Nature* quotes for an estimated male speaker; Light Grey: Proportion of *The Guardian* quotes for an estimated male speaker; Yellow: Proportion of first author articles from an estimated male author in *Springer Nature*; Dark Mustard: Proportion of last author articles from an estimated male author in *Springer Nature*. We observe a larger gender difference between first and last authors in *Springer Nature* articles, however the proportion of predicted male speakers is less than observed in *Nature* research articles. Panel B depicts the proportion of male quotes broken down by article type. In all plots the colored bands represent the 5th and 95th bootstrap quantiles and the point is the mean calculated from 5,000 bootstrap samples.

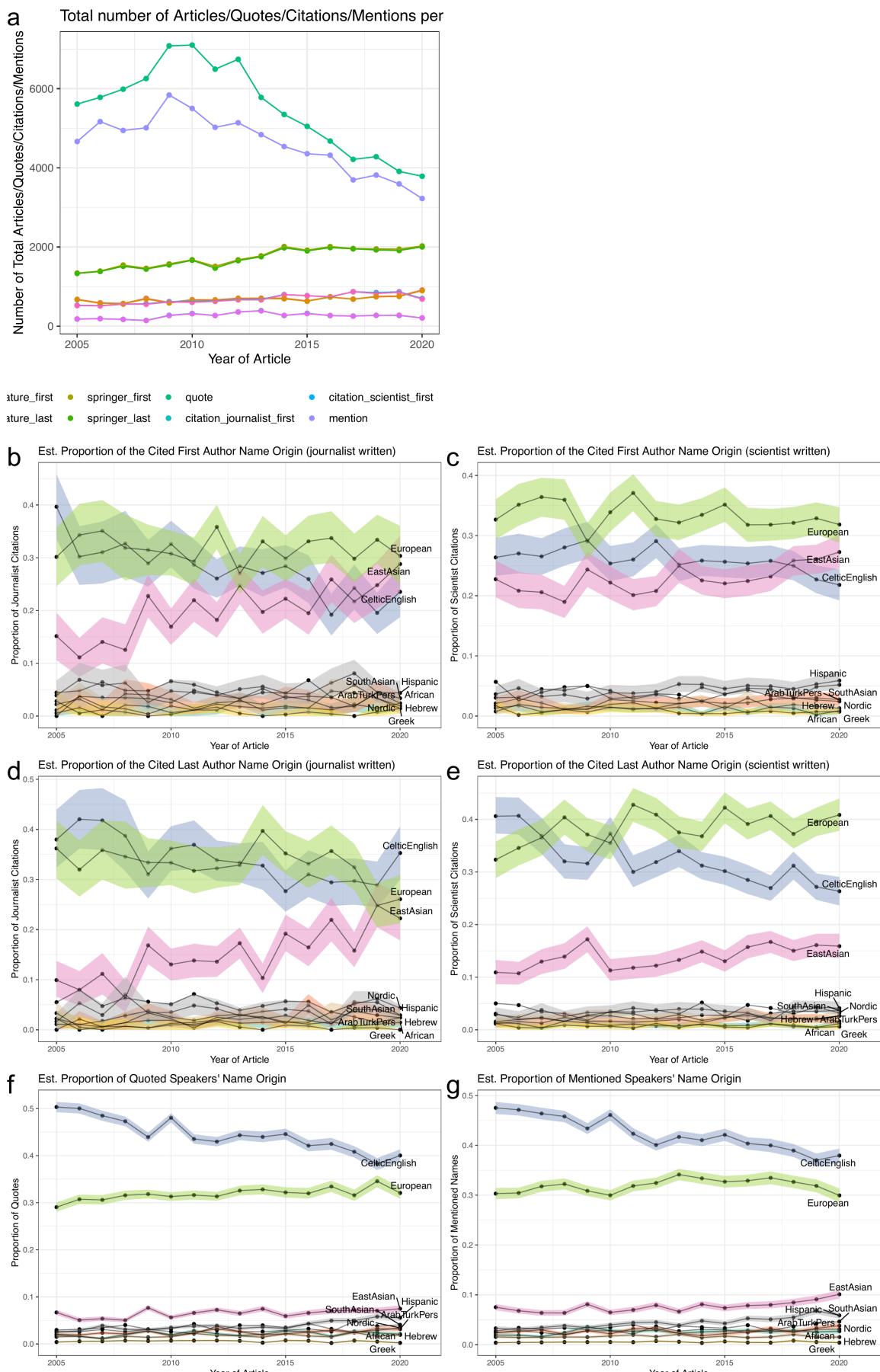


Figure Supplemental 4: Predicted Celtic/English, and European name origins are the highest cited, quoted, and mentioned Panel A, depicts the number of quotes, mentions, citations, or research articles considered in the name origin analysis. Panels B-G depicts the proportion of a name origin in a given dataset, citations in articles written by journalists or writers, quoted speakers or mentions. In all plots the colored bands represent the 5th and 95th bootstrap quantiles and the point is the mean calculated from 5,000 bootstrap samples.

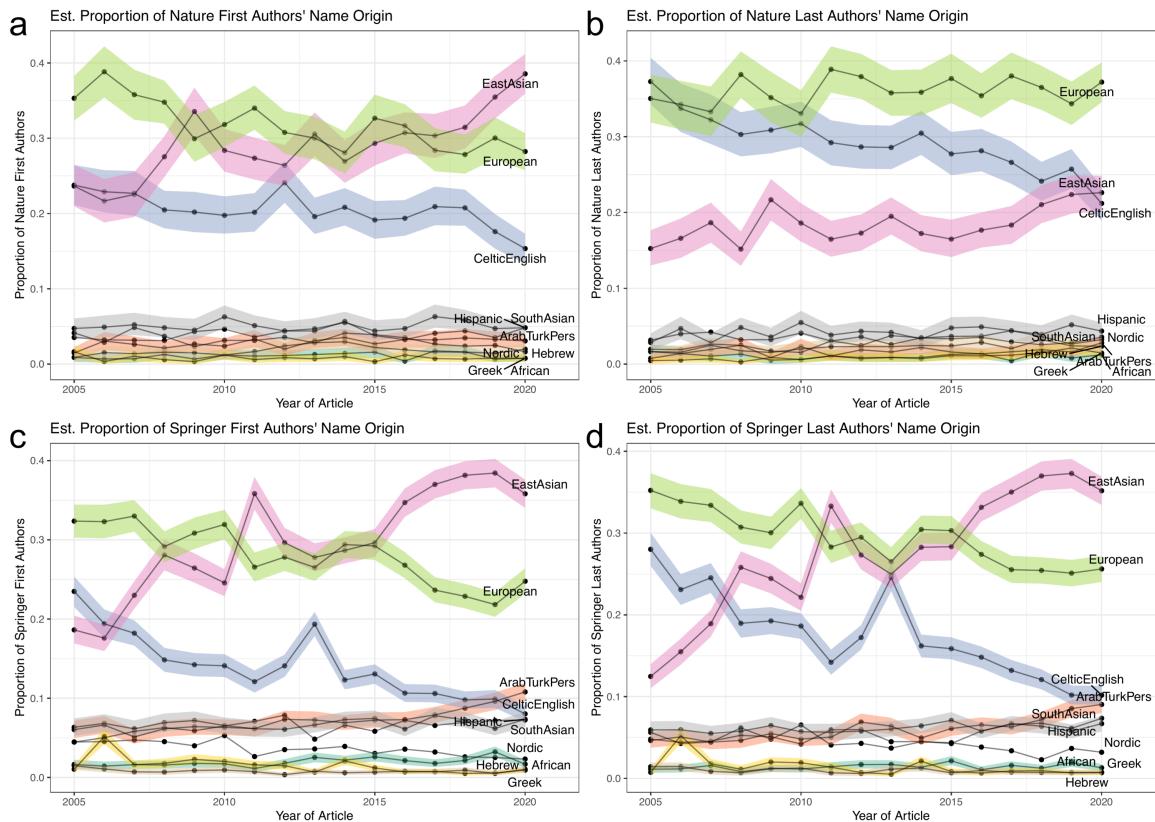


Figure Supplemental 5: Distribution of name origins *Nature* and *Springer Nature* articles Panels A-D depicts the predicted name origins of first and last authors in our background sets. Panel A and B show the predicted name origins of *Nature* first and last authors, respectively. Panel C and D show the predicted name origins of *Springer Nature* first and last authors, respectively.

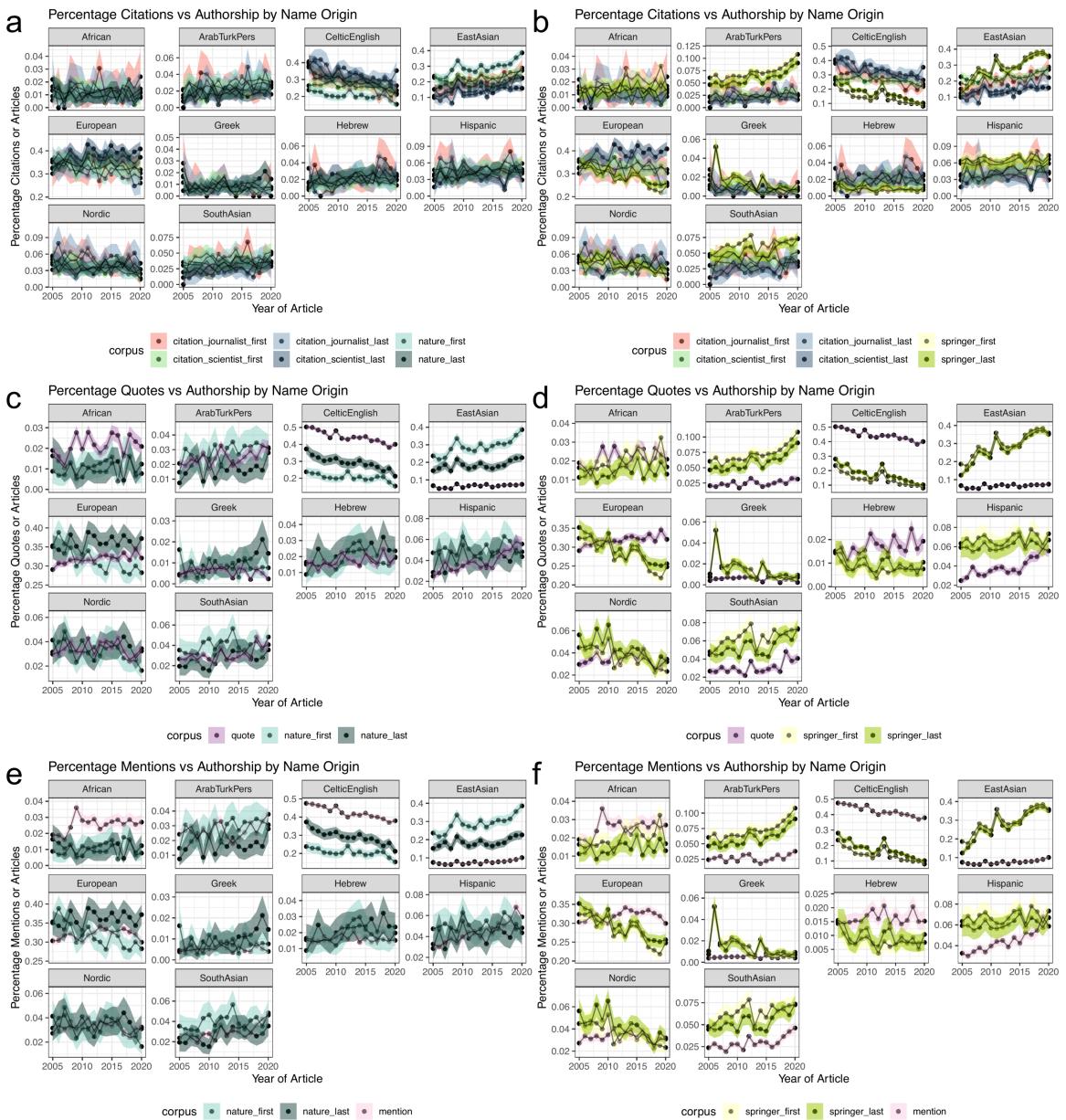


Figure Supplemental 6: Over-representation of predicted Celtic/English and under-representation of East Asian name origins is also found in comparison to *Nature* and *Springer Nature* articles Panels A-F depicts ten plots, each for a possible name origin comparison against a background set. Panel A, C, and E compare the citation (a), quote (c), or mention (e) rate against *Nature* first and last author name origins. Panel B, D, and F compare the citation (a), quote (c), or mention (e) rate against *Springer Nature* first and last author name origins. Panels A and B additionally partition the citation rates by journalist-written articles and scientist-written articles, each further divided into first or last author position. For C-F, only journalist written articles are considered.

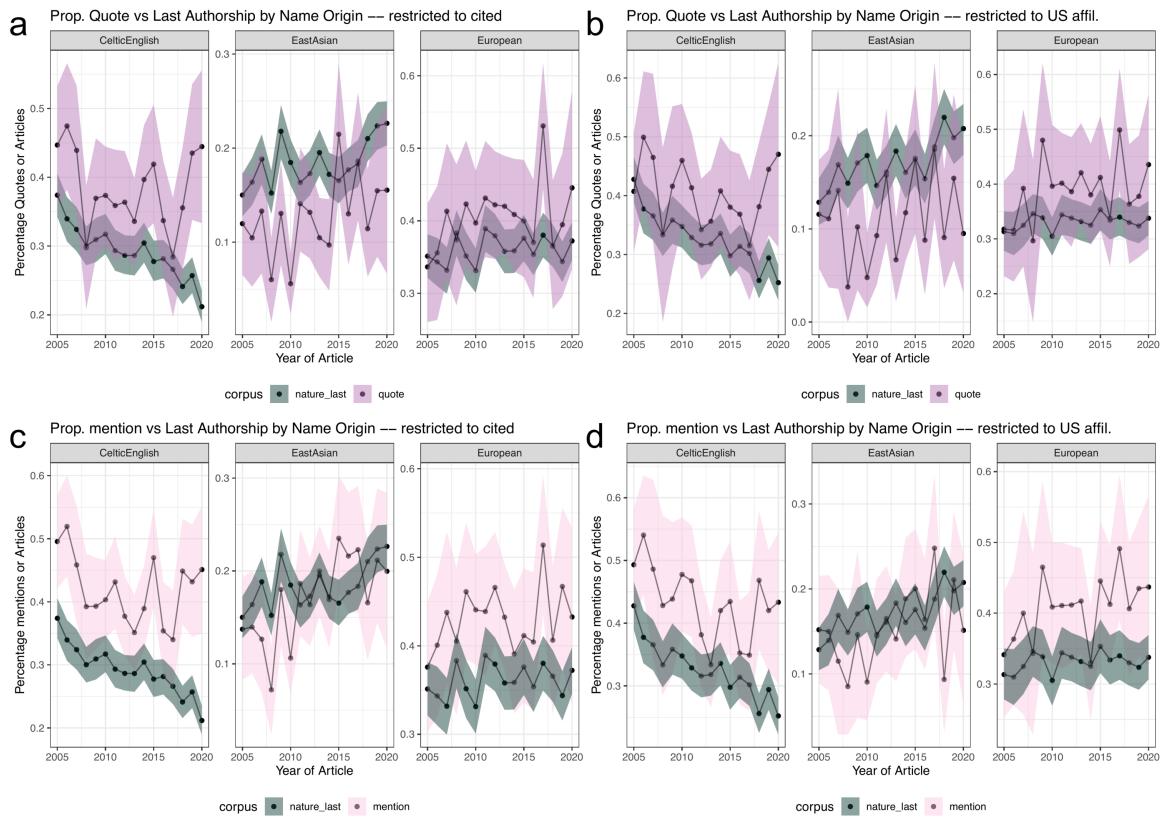


Figure Supplemental 7: Over-representation of predicted Celtic/English and under-representation of East Asian quotes and mentions are reduced when additionally considering citation Panels A-D depicts twelve plots, each for a possible name origin comparison against a background set. Panels A and B compare name origin proportions of quotes from people that were also cited in the same article. Panels C and D compare name origin proportions from mentions of people that were also cited in the same article. In all plots the colored bands represent the 5th and 95th bootstrap quantiles and the point is the mean calculated from 5,000 bootstrap samples.

Table 4: Breakdown of quotes at major processing steps

Processing Step	Frequency
Total Quotes	105457
Quotes with a full name or pronoun associated	96620
Quotes with a gender prediction	96390
Quote with a full name	88535
Quotes with a name origin prediction	100457

Table 5: Breakdown of citations at major processing steps

Writer of Article	Total citations	Total Springer Nature citations	First author citations with a full name	Last author citations with a full name	First author citations with a name origin predictiton	Last author citations with a name origin predictiton
Journalist	15713	5736	4452	4464	4449	4447
Scientist	40707	14597	11276	11170	11276	11152

Table 6: Breakdown of all Springer Nature papers at major processing steps

Processing Step	Frequency
# Springer Nature Articles	38400

Processing Step	Frequency
# First + last authors with a full name in Springer Nature Articles	55370
# First + last authors with a gender prediction in Springer Nature Articles	51686
# First + last authors with a name origin prediction in Springer Nature Articles	55197

Table 7: Breakdown of all Nature papers at major processing steps

Processing Step	Frequency
# Nature Articles	13414
# First + last authors with a full name in Nature Articles	21996
# First + last authors with a gender prediction in Nature Articles	21173
# First + last authors with a name origin prediction in Nature Articles	21996

Table 8: Quoted speaker name origin, by journalist name origin

Journalist Name Origin	African	Arab Turk Pers	Celtic English	East Asian	Europe an	Greek	Hebre w	Hispani c	Nordic	South Asian
CelticEnglish	0.020	0.025	0.484	0.038	0.319	0.006	0.016	0.033	0.035	0.022
EastAsian	0.018	0.017	0.354	0.243	0.250	0.004	0.016	0.026	0.036	0.035
European	0.022	0.023	0.420	0.086	0.326	0.005	0.016	0.043	0.032	0.027

Table 9: Quoted + cited speaker name origin, by journalist name origin %

Journalist Name Origin	African	Arab Turk Pers	Celtic English	East Asian	Europe an	Greek	Hebre w	Hispani c	Nordic	South Asian
CelticEnglish	0.016	0.027	0.368	0.070	0.363	0.008	0.017	0.023	0.083	0.025
EastAsian	0.002	0.077	0.377	0.142	0.167	0.000	0.012	0.133	0.019	0.080
European	0.014	0.028	0.363	0.116	0.352	0.006	0.030	0.026	0.035	0.030

Table 10: Quoted speakers (with US affiliated citation) name origin, by journalist name origin

Journalist Name Origin	African	Arab Turk Pers	Celtic English	East Asian	Europe an	Greek	Hebre w	Hispani c	Nordic	South Asian
CelticEnglish	0.010	0.023	0.378	0.087	0.361	0.010	0.021	0.029	0.056	0.024
EastAsian	0.000	0.066	0.340	0.148	0.209	0.000	0.005	0.148	0.033	0.049
European	0.020	0.030	0.410	0.111	0.300	0.012	0.023	0.019	0.030	0.046

References

1. **Sci-Hub provides access to nearly all scholarly literature**
Daniel S Himmelstein, Ariel Rodriguez Romero, Jacob G Levernier, Thomas Anthony Munro, Stephen Reid McLaughlin, Bastian Greshake Tzovaras, Casey S Greene
eLife (2018-03-01) <https://doi.org/ckcj>
DOI: [10.7554/elife.32822](https://doi.org/10.7554/elife.32822) · PMID: [29424689](#) · PMCID: [PMC5832410](#)
2. **Reproducibility of computational workflows is automated using continuous analysis**
Brett K Beaulieu-Jones, Casey S Greene
Nature biotechnology (2017-04) <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6103790/>
DOI: [10.1038/nbt.3780](https://doi.org/10.1038/nbt.3780) · PMID: [28288103](#) · PMCID: [PMC6103790](#)
3. **Bitcoin for the biological literature.**
Douglas Heaven
Nature (2019-02) <https://www.ncbi.nlm.nih.gov/pubmed/30718888>
DOI: [10.1038/d41586-019-00447-9](https://doi.org/10.1038/d41586-019-00447-9) · PMID: [30718888](#)
4. **Plan S: Accelerating the transition to full and immediate Open Access to scientific publications**
cOAlition S
(2018-09-04) <https://www.wikidata.org/wiki/Q56458321>
5. **Open access**
Peter Suber
MIT Press (2012)
ISBN: 9780262517638
6. **Open collaborative writing with Manubot**
Daniel S Himmelstein, Vincent Rubinetti, David R Slochower, Dongbo Hu, Venkat S Malladi, Casey S Greene, Anthony Gitter
Manubot (2020-05-25) <https://greenelab.github.io/meta-review/>
7. **Opportunities and obstacles for deep learning in biology and medicine**
Travers Ching, Daniel S Himmelstein, Brett K Beaulieu-Jones, Alexandr A Kalinin, Brian T Do, Gregory P Way, Enrico Ferrero, Paul-Michael Agapow, Michael Zietz, Michael M Hoffman, ...
Casey S Greene
Journal of The Royal Society Interface (2018-04) <https://doi.org/gddkhn>
DOI: [10.1098/rsif.2017.0387](https://doi.org/10.1098/rsif.2017.0387) · PMID: [29618526](#) · PMCID: [PMC5938574](#)
8. **Open collaborative writing with Manubot**
Daniel S Himmelstein, Vincent Rubinetti, David R Slochower, Dongbo Hu, Venkat S Malladi, Casey S Greene, Anthony Gitter
PLOS Computational Biology (2019-06-24) <https://doi.org/c7np>
DOI: [10.1371/journal.pcbi.1007128](https://doi.org/10.1371/journal.pcbi.1007128) · PMID: [31233491](#) · PMCID: [PMC6611653](#)
9. **The enduring whiteness of the American media**
Howard W French
The Guardian (2016-05-25) <https://www.theguardian.com/world/2016/may/25/enduring-whiteness-of-american-journalism>
10. **I Analyzed a Year of My Reporting for Gender Bias and This Is What I Found**
Adrienne LaFrance

LadyBits on Medium (2013-09-30) <https://medium.com/ladybits-on-medium/i-analyzed-a-year-of-my-reporting-for-gender-bias-and-this-is-what-i-found-a16c31e1cdf>

11. **I Analyzed a Year of My Reporting for Gender Bias (Again)**

Adrienne LaFrance

The Atlantic (2016-02-17) <https://www.theatlantic.com/technology/archive/2016/02/gender-diversity-journalism/463023/>

12. **I Spent Two Years Trying to Fix the Gender Imbalance in My Stories**

Ed Yong

The Atlantic (2018-02-06) <https://www.theatlantic.com/science/archive/2018/02/i-spent-two-years-trying-to-fix-the-gender-imbalance-in-my-stories/552404/>

13. **A Paper Ceiling**

Eran Shor, Arnout van de Rijt, Alex Miltsov, Vivek Kulkarni, Steven Skiena

American Sociological Review (2015-09-30) <https://doi.org/f7tzps>

DOI: [10.1177/0003122415596999](https://doi.org/10.1177/0003122415596999)

14. **Time Trends in Printed News Coverage of Female Subjects, 1880–2008**

Eran Shor, Arnout van de Rijt, Charles Ward, Aharon Blank-Gomel, Steven Skiena

Journalism Studies (2013-09-12) <https://doi.org/gj3z8b>

DOI: [10.1080/1461670x.2013.834149](https://doi.org/10.1080/1461670x.2013.834149)

15. **Women and news: A long and winding road**

Karen Ross, Cynthia Carter

Media, Culture & Society (2011-11) <https://doi.org/ccxhvz>

DOI: [10.1177/0163443711418272](https://doi.org/10.1177/0163443711418272)

16. **Women Are Seen More than Heard in Online Newspapers**

Sen Jia, Thomas Lansdall-Welfare, Saatviga Sudhahar, Cynthia Carter, Nello Cristianini

PLOS ONE (2016-02-03) <https://doi.org/f8q47g>

DOI: [10.1371/journal.pone.0148434](https://doi.org/10.1371/journal.pone.0148434) · PMID: [26840432](https://pubmed.ncbi.nlm.nih.gov/26840432/) · PMCID: [PMC4740422](https://pubmed.ncbi.nlm.nih.gov/PMC4740422/)

17. **Lack of female sources in NY Times front-page stories highlights need for change**

Poynter

(2013-07-16) <https://www.poynter.org/reporting-editing/2013/lack-of-female-sources-in-new-york-times-stories-spotlights-need-for-change/>

18. **Who Makes the News | GMMP 2015 Reports** <https://whomakesthenews.org/gmmp-2015-reports/>

19. **Women, Minorities, and Persons with Disabilities in Science and Engineering: 2021 | NSF - National Science Foundation** <https://ncses.nsf.gov/pubs/nsf21321/>

20. **Why we need to increase diversity in the immunology research community**

Akiko Iwasaki

Nature Immunology (2019-08-19) <https://doi.org/gkmwwv>

DOI: [10.1038/s41590-019-0470-6](https://doi.org/10.1038/s41590-019-0470-6) · PMID: [31427777](https://pubmed.ncbi.nlm.nih.gov/31427777/)

21. **Men Set Their Own Cites High: Gender and Self-citation across Fields and over Time**

Molly M King, Carl T Bergstrom, Shelley J Correll, Jennifer Jacquet, Jevin D West

Socius: Sociological Research for a Dynamic World (2017-01-01) <https://doi.org/ddzq>

DOI: [10.1177/2378023117738903](https://doi.org/10.1177/2378023117738903)

22. **Bibliometrics: Global gender disparities in science**

Vincent Larivière, Chaoqun Ni, Yves Gingras, Blaise Cronin, Cassidy R Sugimoto

Nature (2013-12) <https://doi.org/gqf>
DOI: [10.1038/504211a](https://doi.org/10.1038/504211a) · PMID: [24350369](https://pubmed.ncbi.nlm.nih.gov/24350369/)

23. **Fund Black scientists**

Kelly R Stevens, Kristyn S Masters, PI Imoukhuede, Karmella A Haynes, Lori A Setton, Elizabeth Cosgriff-Hernandez, Muyinatu A Lediju Bell, Padmini Rangamani, Shelly E Sakiyama-Elbert, Stacey D Finley, ... Omolola Eniola-Adefeso
Cell (2021-02) <https://doi.org/ghvqv5>
DOI: [10.1016/j.cell.2021.01.011](https://doi.org/10.1016/j.cell.2021.01.011) · PMID: [33503447](https://pubmed.ncbi.nlm.nih.gov/33503447/)

24. **NIH peer review: Criterion scores completely account for racial disparities in overall impact scores**

Elena A Erosheva, Sheridan Grant, Mei-Ching Chen, Mark D Lindner, Richard K Nakamura, Carole J Lee
Science Advances (2020-06-05) <https://doi.org/gjnjbz>
DOI: [10.1126/sciadv.aaz4868](https://doi.org/10.1126/sciadv.aaz4868) · PMID: [32537494](https://pubmed.ncbi.nlm.nih.gov/32537494/) · PMCID: [PMC7269672](https://pubmed.ncbi.nlm.nih.gov/PMC7269672/)

25. **Topic choice contributes to the lower rate of NIH awards to African-American/black scientists**

Travis A Hoppe, Aviva Litovitz, Kristine A Willis, Rebecca A Meseroll, Matthew J Perkins, Blan Hutchins, Alison F Davis, Michael S Lauer, Hannah A Valentine, James M Anderson, George M Santangelo
Science Advances (2019-10-04) <https://doi.org/gghp8t>
DOI: [10.1126/sciadv.aaw7238](https://doi.org/10.1126/sciadv.aaw7238) · PMID: [31633016](https://pubmed.ncbi.nlm.nih.gov/31633016/) · PMCID: [PMC6785250](https://pubmed.ncbi.nlm.nih.gov/PMC6785250/)

26. **The Gender Gap in NIH Grant Applications**

Timothy J Ley, Barton H Hamilton
Science (2008-12-05) <https://doi.org/frdj6k>
DOI: [10.1126/science.1165878](https://doi.org/10.1126/science.1165878) · PMID: [19056961](https://pubmed.ncbi.nlm.nih.gov/19056961/)

27. **Race, Ethnicity, and NIH Research Awards**

Donna K Ginther, Walter T Schaffer, Joshua Schnell, Beth Masimore, Faye Liu, Laurel L Haak, Raynard Kington
Science (2011-08-19) <https://doi.org/csf8j8>
DOI: [10.1126/science.1196783](https://doi.org/10.1126/science.1196783) · PMID: [21852498](https://pubmed.ncbi.nlm.nih.gov/21852498/) · PMCID: [PMC3412416](https://pubmed.ncbi.nlm.nih.gov/PMC3412416/)

28. **Including Diverse Voices in Science Stories**

Christina Selby
The Open Notebook (2016-08-23) <https://www.theopennotebook.com/2016/08/23/including-diverse-voices-in-science-stories/>

29. **gage. Discover Brilliance** <https://gage.500womenscientists.org/>

30. **WMC SheSource - Women's Media Center** <https://womensmediacenter.com/shesource>

31. **Scrapy | A Fast and Powerful Scraping and Web Crawling Framework** <https://scrapy.org/>

32. **The Stanford CoreNLP Natural Language Processing Toolkit**

Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, David McClosky
Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations (2014) <https://doi.org/gf3xhp>
DOI: [10.3115/v1/p14-5010](https://doi.org/10.3115/v1/p14-5010)

33. **humanifomat: A Parser for Human Names**

Oliver Keyes

(2016-04-24) <https://CRAN.R-project.org/package=humaniformat>

34. **Gender Prediction Methods Based on First Names with genderizeR**
Kamil Wais
The R Journal (2016) <https://doi.org/gf4zqx>
DOI: [10.32614/rj-2016-002](https://doi.org/10.32614/rj-2016-002)
35. **Genderize.io | Determine the gender of a name** <https://genderize.io/>
36. **Analysis of ISCB honorees and keynotes reveals disparities**
Trang T Le, Daniel S Himmelstein, Ariel A Hippen Anderson, Matthew R Gazzara, Casey S Greene
Cold Spring Harbor Laboratory (2020-04-14) <https://doi.org/ggr64p>
DOI: [10.1101/2020.04.14.927251](https://doi.org/10.1101/2020.04.14.927251)
37. **Nationality Classification Using Name Embeddings**
Junting Ye, Shuchu Han, Yifan Hu, Baris Coskun, Meizhu Liu, Hong Qin, Steven Skiena
Proceedings of the 2017 ACM on Conference on Information and Knowledge Management
(2017-11-06) <https://doi.org/ggjc78>
DOI: [10.1145/3132847.3133008](https://doi.org/10.1145/3132847.3133008)
38. **OpenStreetMap**
OpenStreetMap
<https://www.openstreetmap.org/>
39. **Comparison and benchmark of name-to-gender inference services**
Lucía Santamaría, Helena Mihaljević
PeerJ Computer Science (2018-07-16) <https://doi.org/ggpbn6>
DOI: [10.7717/peerj-cs.156](https://doi.org/10.7717/peerj-cs.156) · PMID: [33816809](#) · PMCID: [PMC7924484](#)
40. **Using genderize.io to infer the gender of first names: how to improve the accuracy of the inference**
Paul Sebo
Journal of the Medical Library Association (2021-11-22) <https://doi.org/gn94vp>
DOI: [10.5195/jmla.2021.1252](https://doi.org/10.5195/jmla.2021.1252) · PMID: [34858090](#) · PMCID: [PMC8608220](#)
41. **Underrepresentation of Asian awardees of United States biomedical research prizes**
Yuh Nung Jan
Cell (2022-02) <https://doi.org/gpgvs2>
DOI: [10.1016/j.cell.2022.01.004](https://doi.org/10.1016/j.cell.2022.01.004) · PMID: [35120660](#)
42. **Racial and ethnic imbalance in neuroscience reference lists and intersections with gender**
Maxwell A Bertolero, Jordan D Dworkin, Sophia U David, Claudia López Lloreda, Pragya Srivastava, Jennifer Stiso, Dale Zhou, Kafui Dzirasa, Damien A Fair, Antonia N Kaczkurkin, ... Danielle S Bassett
Cold Spring Harbor Laboratory (2020-10-12) <https://doi.org/gj7mdc>
DOI: [10.1101/2020.10.12.336230](https://doi.org/10.1101/2020.10.12.336230)