**Figure 1:** Sentence reuse across two articles on "self-driving cars". Our pipeline detected reuse even across edits, e.g.: *"...Musk's own **statements issuing** lofty goals for production cars or turning a sustained profit starting this quarter"* vs *"...Musk's own **pronouncements such as** lofty goals for production cars or turning a sustained profit starting this quarter **that might be beyond reach"** detected as variants.

# News Provenance: Revealing News Text Reuse at Web-Scale in an Augmented News Search Experience

**Nathan Evans**
Microsoft Research
Silverdale, WA, USA
naevans@microsoft.com

**Jonathan Larson**
Microsoft Research
Silverdale, WA, USA
jolarso@microsoft.com

**Darren Edge**
Microsoft Research
Cambridge, UK
daedge@microsoft.com

**Christopher White**
Microsoft Research
Redmond, WA, USA
chwh@microsoft.com

## Abstract

The media industry has a practice of reusing news content, which may be a surprise to news consumers. Whether by agreement or plagiarism, a lack of explicit citations makes it difficult to understand where news comes from and how it spreads. We reveal *news provenance* by reconstructing the history of near-duplicate news in the web index – identifying the origins of republished content and the impact of original content. By aggregating provenance information and presenting it as part of news search results, users may be able to make more informed decisions about which articles to read and which publishers to trust. We report on early analysis and user feedback, highlighting the critical tension between the desire for media transparency and the risks of disrupting an already fragile ecosystem.

## Author Keywords

news media; news search; media transparency; provenance analysis; provenance visualization

## CCS Concepts

•**Human-centered computing** → *Activity centered design; Visual analytics;* •**Information systems** → *Near-duplicate and plagiarism detection; Data provenance; Presentation of retrieval results;*
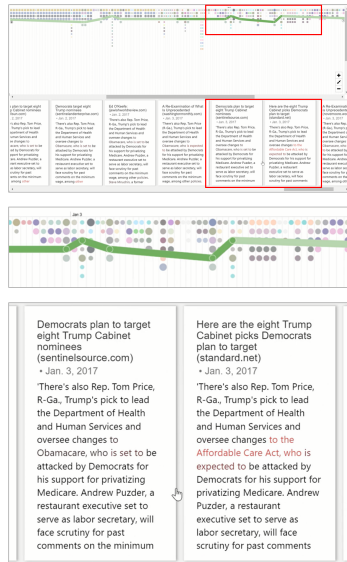
**Figure 2:** News Threads visual for Power BI (not released). Articles are shown as vertical dot plots of individual paragraphs and ordered by publication time. Paragraphs are clustered by similarity, with dots colored and sized based on cluster identity and republication count. Selecting a paragraph's dot reveals its provenance "thread" and lists clustered paragraphs for comparison. TF-IDF is applied at the trigram level to highlight (in red) text that is distinctive within a cluster, e.g., *"...changes to* ***Obamacare*** *who is* ***set*** *to be attacked..."* vs *"...changes to* ***the Affordable Care Act*** *who is* ***expected*** *to be attacked..."*.

## Introduction

Since the first newspapers were printed in the seventeenth century, following the news has become a central activity of everyday life. While the shift to digital distribution has increased both the production and consumption of news, the pressures of the 24-hour news cycle combined with shrinking revenues and resources has led to a "crisis in journalism" and a "publish first, verify later" culture [17].

The consequences of this crisis are many, from news plagiarism [16] to state-sponsored disinformation campaigns [1]. We are interested in a third phenomenon largely ignored in the public debate, however: the practice of journalistic text reuse. Earlier work contrasted the "benign and acceptable" reuse of news agency text ("copy") from the more harmful and unacceptable practice of plagiarism [5]. Today, however, content can easily be copied from anywhere by anyone, without agreement and with little chance of repercussions. The resulting potential for mass republication creates a multitude of threats – an economic threat to news publishers, a political threat to democratic processes, and an epistemological threat to news consumers who no longer know what to believe. For publishers and consumers alike, it has never been more important to understand *news provenance* – where news content has come from and how it has been modified and republished over time.

This case study reports on our early efforts to reveal such news provenance based on automated analysis of news article text – identifying original and duplicated sentences in ways that allow us to create a rich graph of content reuse across both articles and publishers. Our work takes advantage of the fact that Microsoft is one of only several organizations with a broad and near real-time index of online news, combined with a search engine and news portal through which many millions of users find and consume news every day. It also builds on our prior work developing text analytics capabilities for Microsoft Power BI [7], where we observed first-hand the difficulty of making sense of a news ecosystem distorted by republication (and designed a Power BI visual to explore patterns of reuse – Figure 2).

## Related Work

News text reuse has been common practice since the era of print publication. The 1999 METER project [5] presents the most in-depth study of the phenomenon, contributing a corpus of 772 Press Association articles and 944 articles from 9 UK daily newspapers, with both article and sentence level reuse annotated by domain experts. The project reported that 78% of news articles were derived from news agency sources (46% partially, 32% wholly) and introduced a variety of techniques (n-gram overlap, greedy string tiling, and sentence alignment) for classifying the extent of derivation.

The concept of data provenance is similarly well known, gaining prominence with the development of the semantic web and motivated by the observed reuse of news text [10]. The resulting data model, PROV-DM, has since been used to describe the provenance of news stories as inferred from the named entities mentioned within them (evaluated on a week of 410 publications by a single publisher) [6].

For the problem of identifying shared text rather than shared entities, Locality-Sensitive Hashing (LSH) of individual sentences has been proposed as a logical but expensive approach [15], with document fingerprinting based on the Discrete Cosine Transform (DCT) preferred for detecting 'local' text reuse across a newswire collection (Table 1).

Analyzing news provenance on the web presents additional challenges in terms of document quantity, quality, and variety, as well as the need to retrieve the appropriate document set for analysis. One approach to web text

| Reuse type | Prevalence |
|------------|------------|
| *Most–Most* | 8% |
| *Most–Much* | 10% |
| *Most–Some* | 10% |
| *Much–Much* | 17% |
| *Much–Some* | 27% |
| *Some–Some* | 28% |

**Table 1:** Text reuse types and their relative prevalence in a TREC newswire collection of 758,224 articles by the Associated Press, Wall Street Journal, Financial Times, and Los Angeles Times, 1987–1994 [15]. Reuse is defined in terms of mutual containment over article pairs, with *Most* as 80-100% containment, *Much* as 50-80%, and *Some* as 10-50%. 51% of articles overlapped with at least one other article.

provenance [3] is to systematically construct, evaluate, and rank the results of multiple search queries comprising structured subsets of the target text based on its noun phrases and named entities. While this approach can leverage the comprehensive web index of the underlying search engine, such indices are not designed with text provenance in mind and cannot be queried at the scale required to reconstruct provenance graphs for large document collections.

The challenge of news provenance on the web is exacerbated by the growing number of publishers arising from the low barriers to entry. Recent work [11] describes how mass news republication by different sources can distort perceptions of story significance, or develop a sense of credibility for an otherwise unreliable source which may combine legitimate, unattributed articles with alternative, fake, or conspiracy-oriented news. The authors study this phenomenon by analyzing verbatim text reuse across 54k political news articles in 2017. Of 92 publishers, 67 were found to republish at least one article, with reuse occurring more frequently among publishers with similar audiences (e.g., for mainstream vs alternative vs hyperpartisan news).

A study of user evaluation strategies for news articles appearing in the Facebook news feed [8] has shown that the source of an article is a critical component for evaluating its trustworthiness, yet readers rarely question whether the listed source is accurate and legitimate. One successful approach to representing credibility in the context of search results [14] is to augment results with visual badges, ratings, or rankings of attributes that would be difficult or impossible for users to assess directly.

In summary, prior work has studied news provenance and search result augmentation, but never in combination as a means of supporting the real-world activities of both news publishers and consumers.

## Implementation of News Provenance

Our prototyping of news provenance capabilities has focused on the iterative and incremental development of data pipelines supporting differentiated user experiences.

*Common Data Pipeline*

Our technical approach is based on the NLP technique of *w-shingling* for computing text similarity – extracting fixed-length n-grams from each comparable text unit and taking the Jaccard similarity (intersection divided by union) of these token sets. We use a Locality-Sensitive Hashing (LSH) index of n-gram token sets represented as MinHash values, which provides sub-linear estimation of Jaccard similarity for any given text query. For a given similarity threshold, our pipeline will cluster similar text units and sort them by the publication time of the containing article, thus reconstructing the provenance and republication of all observed text units. This provenance information can then be aggregated to the article and publisher levels to give a complete picture of text reuse within the news ecosystem.

By first applying our pipeline at the document level, we can identify independent clusters of documents within which sentence clustering is computationally feasible (since clustering requires pairwise comparisons that grow much more rapidly for sentences than documents). Generalizing this idea, our pipeline supports efficient hierarchical clustering of arbitrary text units. It first populates a sparse matrix with pairwise text distances when the Jaccard similarity exceeds a certain threshold, before applying DBSCAN clustering in a hierarchical fashion. At each iteration, the espilon parameter is reduced to progressivly decrease the variation within each identified cluster, and clustering only takes place within previously identified clusters. The result is a hierarchical clustering moving from all text units at the root to duplicate (or near-duplicate) text units at the leaves.

| Reuse type | Prevalence |
|------------|-----------|
| *Most–Most* | 46% |
| *Most–Much* | 7% |
| *Most–Some* | 4% |
| *Much–Much* | 7% |
| *Much–Some* | 5% |
| *Some–Some* | 31% |

**Table 2:** Text reuse types and their relative prevalence as analyzed by our news provenance pipeline over a corpus of 45,300 news articles matching "Kamala Harris" from May–July 2019. 97% of articles were derived ($\geq$10% sentence similarity with other articles), as were 97% of sentences ($\geq$50% trigram similarity with other sentences). Reuse types were classified using the scheme from [15], aggregated over all pairwise article comparisons, and normalized to 100%. In contrast to prior results from the era of print publication (Table 1), we can see that the newly dominant form of reuse is one in which entire articles are copied directly with only minor edits and additions.
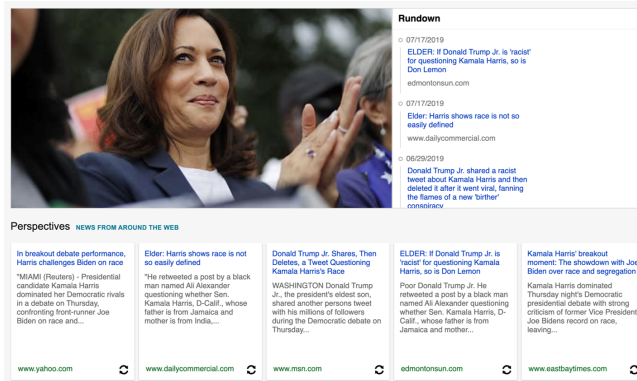


**Figure 3:** Consumer experience for the news query "Kamala Harris", showing a timeline of published articles (right) and a card carousel for clusters of similar articles (bottom).

*Consumer Experience*
The Bing News Spotlight feature [18] provides a model for giving news consumers an overview of major developing news stories, including a topic timeline and presentation of different "perspectives". We adopted the same model for the design of our consumer experience (Figure 3). Clicking the icon in the lower-right corner of a perspective card reveals more information about the article cluster (Figure 4), while opening the article reveals sentence-level provenance (Figure 5). Selecting two articles shows a side-by-side view of fuzzy sentence matches and their differences (Figure 1).

To recreate the "perspectives" view based on clusters of similar articles at the text level, we configure our pipeline to operate on document unigrams (i.e., focusing on the presence but not order of words) and use our adaptive clustering method to identify clusters of related articles. Sentences are then clustered as near-duplicates when their trigram sets have a Jaccard similarity of 0.5 or greater.
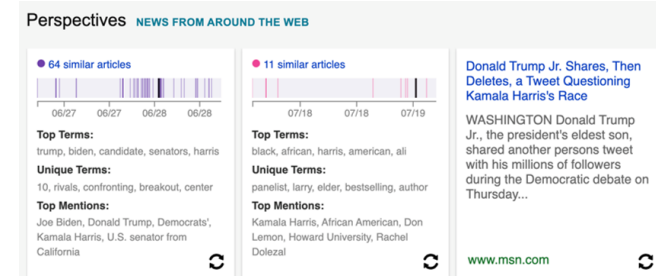


**Figure 4:** Cards can be flipped to reveal more details about the article cluster. Selecting "N similar articles" drills down to the individual articles of that cluster.

*Publisher Experience*
Publishers deeply understand the ecosystem of news wires, syndication, and editorial decision-making in a way that most consumers do not. Our publisher experience therefore includes support for verifying sources, understanding content propagation, and revealing both gaps in reporting coverage and patterns in publication timing. The interface shows publication timelines for content variants, navigable document clusters, and a treemap of news domains encoding both popularity and share of original content (Figure 6).

*Analyst Experience*
The Power BI visual we created to explore "news threads" (Figure 2) confirmed what we had informally observed in our prior development of news analytics applications [7] – that near-verbatim duplication of entire paragraphs was widespread across the results of news search queries. We observed likely cases of automated syndication (sequences of identical dot patterns), rolling updates to breaking news (dot patterns shifted downwards over time as new paragraphs were inserted at the top), remixing of content from multiple distinct sources, and edits both benign and slanted.

## Quote variations

*"you know"* harris said turning to former vice president joe biden *"there was..."*

...harris informed him that *"there was..."*

*"...every day"* harris 54 told biden.

*"...every day"* harris whose father is black and mother is of indian descent said.

*"...every day"* then the killer line *"that little girl was me."*

*"...was me"* harris a california democrat said.

*"...was me"* harris said at times looking directly at biden as he gazed elsewhere.

*"...was me"* the moment apparently planned ... led to harris being unofficially declared the debate winner.

**Table 3:** Example sentence variants reporting the Kamala Harris quote *"and you know there was a little girl in california who was part of the second class to integrate her public schools and she was bused to school every day and that little girl was me"*. Different spellings, insertions/deletions, and paraphrasing also cause variation.



**Figure 5:** Selecting the article title shows each sentence of the article annotated with the article where it was first observed. Selecting a publisher highlights all associated derivation links.

To support this analyst experience, we configure our data pipeline to operate on paragraph-level text units and output the data tables needed to drive this visual within Power BI.

## Example Analysis

We have analyzed a wide range of news topics including healthcare, self-driving cars, and eSports, as well as more controversial topics such as gun control. Our goal was to understand how news is copied and modified across syndication outlets, local affiliates, and content farms over time, as well as to evaluate and refine our provenance pipeline and the experiences built upon it. Here we describe one such analysis of a controversial news cycle about Kamala Harris – a US Senator whose race and eligibility as a presidential candidate were called into question following the launch of her 2020 US election campaign.
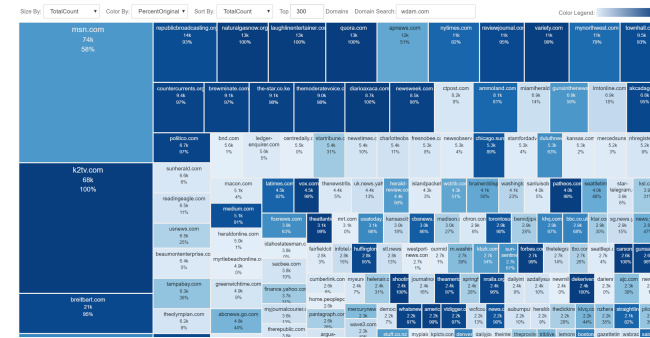


**Figure 6:** Domain treemap from publisher experience. Cell size and color represent article count and original content respectively for the topic "gun control".

We indexed a corpus of 45,300 news articles containing "Kamala Harris" from the period May–July 2019, reliably tracing content provenance back to a widely-copied opinion piece by Larry Elder (Figure 5). Analysis of reuse types using the method of [15] is shown in Table 2. We observe dominance of the *Most–Most* reuse type across the majority of our analyses, confirming a substantial change in news text reuse since the shift to online distribution.

Overall, we found 254,094 unique sentences from 1,522,340 sentence-to-document mappings, giving an average of 34 sentences per document. Clustering these sentences using a trigram similarity threshold of 0.5 resulted in 212,714 sentence clusters, with a mean of 7.2 *sentence publications* per cluster (median 1, maximum 3003). The mean number of *sentence variants* within each cluster was 1.2 (median 1, maximum 105 – see Table 3 for analysis of this cluster).
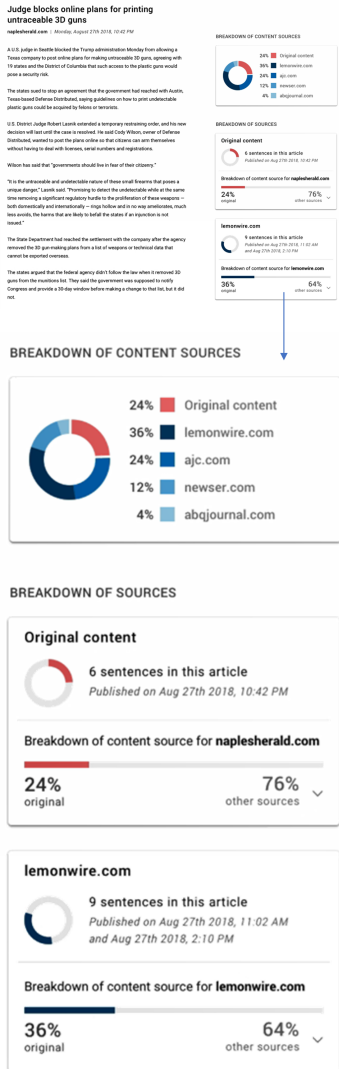
**Figure 7:** Concept video image.

## Concept Testing

We performed both consumer surveys and publisher interviews to understand attitudes towards news provenance. Both used a concept video walk-through of a representative interface showing a news reading experience (Figure 7).

*Consumer Survey*

We recruited 40 participants to complete an online survey (20 male/female, average age 33). All were online news consumers with 78% checking the news at least daily. Participants were assigned to four conditions in a between-subjects design comparing publisher type (local vs national) and degree of originality (mostly original vs mostly sourced). A 7-point scale from extremely low to extremely high was used to understand how these factors may influence user confidence in article content.

We found that while *original local* content increased confidence (9/10 high with 5/10 extremely high), *sourced national* content decreased confidence (4/10 low with 4/10 slightly low). The majority of participants rated the interface as easy to use (35/40 easy with 21/40 extremely easy).

Participants were also asked to rank the concepts of authority, credibility, trustworthiness, and originality in order of importance. 19/40 each ranked trustworthiness and credibility as most important, with 26/40 participants ranking originality as least important. Multiple participants reported that they were indifferent to originality so long as other conditions were met, e.g., "As long as the source is accurate and truthful and well-recognized as credible" and "News is news. If it is accurate, I am less concerned with originality". Some noted how content reuse could increase credibility, e.g., "It makes me think they are more credible if they are willing to pull from other sources. It seems more balanced", while others saw reuse as beneficial only if it came from a trusted source or was accompanied by commentary on

the reused text. Concerns about representing news provenance included accuracy (e.g., "how do you confirm the actual original source of the information?"), bias (e.g., "it might be possible ... to take a side and make people think that some articles are actually more original than others"), misinterpretation (e.g., "just because something is 100% unoriginal, that isn't bad"), and misuse (e.g., "I'd be afraid that this would be used to discredit legitimate websites").

*Publisher Interviews*

We interviewed representatives of eight major news publishers in the United States for critical feedback on the concept of news provenance. All supported the use of provenance interfaces as newsroom tools for research and reporting (e.g., "News rooms would love to see it"), noting how they could be used to identify original sources, showcase information from trusted sources, support proper attribution, detect plagiarism, track stories across networks, observe how different publishers use different sources, and find sources that are cited both widely and rarely.

As a consumer-facing capability, one broad view was that "publishers should have nothing to fear" because "more transparency and [content] tracking is a good thing" that would "incentivize original content" in ways both "good for user and original news org" . Reported benefits to consumers included access to a greater diversity of sources in search results that could also reduce redundancy.

At the same time, publishers raised concerns about possible unintended consequences of such capabilities (e.g., source devaluation and effects on search result rankings) and the difficulties of both calculating and interpreting aggregate measures of original content. "Publishers could get upset" if they felt they were unfairly represented, while "users would be shocked at [the level of] copying". Four publishers also stated that they would like to see an evalua-

tion of source quality included as part of news provenance, although the question of who gets to judge quality and how are open questions that are significantly more contentious than any measures of original versus sourced content.

## Discussion and Conclusion

This case study brings together diverse strands of work and applies them in ways that enrich our understanding of modern day news. It also suggests multiple directions for future work. For example, the analysis of text reuse is not limited to news content – it can be applied to other texts just as similar methods have been applied in the diverse areas of ancient literature, [13], state legislature [4], and Wikipedia articles [2]. The concept of news provenance is also not limited to news text, e.g., the New York Times has announced a blockchain project to track multimedia provenance [12]. Finally, Microsoft is not the only company analyzing originality – Google has also said it will be elevating the ranking of original content in its search results [9].

The results of our news provenance concept testing suggest that while news consumers may prioritize credibility over originality in abstract terms, the explicit representation of original versus sourced content has a concrete impact on consumers' evaluation of the same underlying news text. Great care must therefore be taken when deciding if and how to communicate such high-level provenance statistics. Media transparency might ultimately be best served by using inferred provenance to encourage publishers to provide more complete attribution for their articles, as a "source checking" counterpart to fact checking. Interfaces might also simply refer to the "first observations" of individual text units, rather than report the proportion of "original content" in ways that might be misinterpreted or manipulated. On this final point, while full disclosure of provenance algorithms could reduce misinterpretation, it could also create

confusion and increase the risk of adversarial attacks. All of these factors need to be considered when designing future mechanisms for news provenance.

## Acknowledgements

## REFERENCES

[1] Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives* 31, 2 (2017), 211–36.

[2] Milad Alshomary, Michael Völske, Tristan Licht, Henning Wachsmuth, Benno Stein, Matthias Hagen, and Martin Potthast. 2019. Wikipedia Text Reuse: Within and Without. In *European Conference on Information Retrieval*. Springer, 747–754.

[3] Michael Bendersky and W Bruce Croft. 2009. Finding text reuse on the web. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining*. ACM, 262–271.

[4] Matthew Burgess, Eugenia Giraudy, Julian Katz-Samuels, Joe Walsh, Derek Willis, Lauren Haynes, and Rayid Ghani. 2016. The Legislative Influence Detector: Finding Text Reuse in State Legislation. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 57–66.

[5] Paul Clough, Robert Gaizauskas, Scott SL Piao, and Yorick Wilks. 2002. Meter: Measuring text reuse. In *Proceedings of the 40th Annual Meeting on*

*Association for Computational Linguistics*. Association for Computational Linguistics, 152–159.

[6] Tom De Nies, Sam Coppens, Davy Van Deursen, Erik Mannens, and Rik Van de Walle. 2012. Automatic discovery of high-level provenance using semantic similarity. In *International Provenance and Annotation Workshop*. Springer, 97–110.

[7] Darren Edge, Jonathan Larson, and Christopher White. 2018. Bringing AI to BI: Enabling Visual Analytics of Unstructured Data in a Modern Business Intelligence Platform. In *Extended Abstracts of the 2018 CHI Conference on Human Factors in Computing Systems (CHI EA '18)*. ACM, New York, NY, USA, Article CS02, 9 pages. DOI: http://dx.doi.org/10.1145/3170427.3174367

[8] Martin Flintham, Christian Karner, Khaled Bachour, Helen Creswick, Neha Gupta, and Stuart Moran. 2018. Falling for Fake News: Investigating the Consumption of News via Social Media. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. ACM, New York, NY, USA, Article 376, 10 pages. DOI: http://dx.doi.org/10.1145/3173574.3173950

[9] Richard Gingras. 2019. Elevating original reporting in Search. https://www.blog.google/products/search/original-reporting/. (2019). Published Sep. 12, 2019. Accessed Oct. 1, 2019.

[10] Paul Groth, Yolanda Gil, James Cheney, and Simon Miles. 2012. Requirements for provenance on the web. *International Journal of Digital Curation* 7, 1 (2012), 39–56.

[11] Benjamin D Horne and Sibel Adali. 2018. An Exploration of Verbatim Content Republishing by News Producers. *arXiv preprint arXiv:1805.05939* (2018).

[12] Sasha Koren. 2019. Introducing the News Provenance Project. https://open.nytimes.com/introducing-the-news-provenance-project-723dbaf07c44. (2019). Published Jul. 23, 2019. Accessed Oct. 1, 2019.

[13] John Lee. 2007. A computational model of text reuse in ancient literary texts. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*. 472–479.

[14] Julia Schwarz and Meredith Morris. 2011. Augmenting web pages and search results to support credibility assessment. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 1245–1254.

[15] Jangwon Seo and W Bruce Croft. 2008. Local text reuse detection. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 571–578.

[16] Rui Sousa-Silva. 2015. Reporter fired for plagiarism: a forensic linguistic analysis of news plagiarism. (2015).

[17] Kate Starbird, Dharma Dailey, Owla Mohamed, Gina Lee, and Emma S. Spiro. 2018. Engage Early, Correct More: How Journalists Participate in False Rumors Online During Crisis Events. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. ACM, New York, NY, USA, 105:1–105:12. DOI: http://dx.doi.org/10.1145/3173574.3173679

[18] Bing Team. 2018. Bing helps you learn more about the news in less time. https://blogs.bing.com/search/2018-08/Bing-helps-you-learn-more-about-the-news-in-less-time. (2018). Published Aug. 27, 2018. Accessed Oct. 1, 2019.