

# News Déjà Vu: Connecting Past and Present with Semantic Search

Brevin Franklin, Emily Silcock, Abhishek Arora, Tom Bryan, and Melissa Dell\*

Harvard University, Cambridge, MA, USA. Authors contributed equally.

\*Corresponding author: melissadell@fas.harvard.edu.

## Abstract

Social scientists and the general public often analyze contemporary events by drawing parallels with the past, a process complicated by the vast, noisy, and unstructured nature of historical texts. For example, hundreds of millions of page scans from historical newspapers have been noisily transcribed. Traditional sparse methods for searching for relevant material in these vast corpora, *e.g.*, with keywords, can be brittle given complex vocabularies and OCR noise. This study introduces News Déjà Vu, a novel semantic search tool that leverages transformer large language models and a bi-encoder approach to identify historical news articles that are most similar to modern news queries. News Déjà Vu first recognizes and masks entities, in order to focus on broader parallels rather than the specific named entities being discussed. Then, a contrastively trained, lightweight bi-encoder retrieves historical articles that are most similar semantically to a modern query, illustrating how phenomena that might seem unique to the present have varied historical precedents. Aimed at social scientists, the user-friendly News Déjà Vu package is designed to be accessible for those who lack extensive familiarity with deep learning. It works with large text datasets, and we show how it can be deployed to a massive scale corpus of historical, open-source news articles. We curate some examples on [newsdejavu.github.io](https://newsdejavu.github.io). While human expertise remains important for drawing deeper insights, News Déjà Vu provides a powerful tool for exploring parallels in how people have perceived past and present.

"Those who cannot remember the past are condemned to repeat it." – George Santayana's *The Life of Reason*

## 1 Introduction

Social scientists, and the public more generally, often seek to place the present in perspective by reflecting upon parallels with the past. Finding these

commonalities, however, can be a labor-intensive and challenging process. Vast troves of historical texts have been preserved, but they are often held in unstructured, uncataloged, massive-scale databases. For example, hundreds of millions of pages from historical newspapers have been digitized and are available online through both open-source and proprietary collections. Keyword searches are often used to extract relevant documents from these massive corpora. However, as language is complex and OCR noise is rampant, sparse methods can be extremely brittle.

Transformer large language models offer a powerful tool for retrieving source material from the past that can contextualize the present. This study trains a novel semantic search model, News Déjà Vu, to query which historical news articles are most semantically similar to a modern news article query. Figure 1 shows the model architecture at inference time. Named entities are first detected and masked out, using a named entity recognition model that we tuned for noisy, historical texts. This allows the model to focus on the generalities of the story, rather than specific names of people, locations, or organizations. Then, we use a contrastively trained bi-encoder model to retrieve the modern article's nearest neighbor(s) from a massive-scale database of historical texts.

The News Déjà Vu package allows social scientists to deploy News Déjà Vu on their own queries. It has a CC-BY license and can be used with any appropriately formatted text dataset. It is designed to be user-friendly and intuitive to social scientists, who often lack knowledge of deep learning frameworks. This study provides code snippets showing how it can be used seamlessly with American Stories, a Hugging Face dataset containing over 430 million historical public domain newspaper article texts (Dell et al., 2023). Interested users can query content from a sampling of states in American Stories with modern articles using our

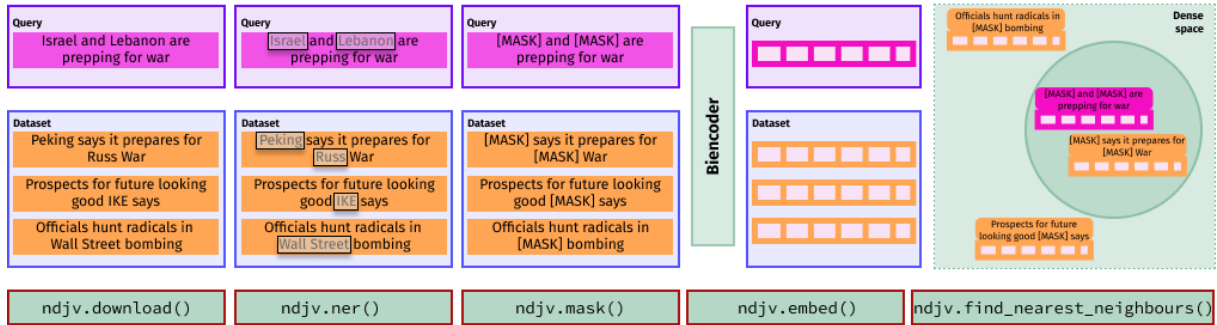


Figure 1: News Déjà Vu architecture at inference time.

HuggingFace demo.<sup>1</sup> We also maintain a website ([newsdejavu.github.io](https://newsdejavu.github.io)), where randomly selected modern news articles - as well as special editions on hand-curated topics of interest - are paired with their retrieved historical neighbors. We built the News Déjà Vu package after many website readers requested that we create a tool that they could use to query their own texts of interest.

News Déjà Vu retrieves articles that use similar semantics. Of course, events or phenomena that are at their core very different may be described in a similar way by the historical and modern news media. This phenomenon is also likely to be of considerable interest to social scientists, but we do caveat that deeper historical knowledge is required to place parallels in appropriate context.

News Déjà Vu currently supports English. In the future, it could be relatively straightforward to create a multilingual model by starting with multilingual Sentence BERT weights and tuning on machine-translated allsides data or other data sources (*e.g.*, [Chen et al. \(2022\)](#)).

The rest of this study is organized as follows. Section 2 discusses the relevant literature, and Section 3 describes the model architecture and training. Section 4 introduces the News Déjà Vu package.

## 2 Related Literature

There is a large literature on semantic similarity. Most large scale datasets in this space are constructed from web texts. The Massive Text Embedding Benchmark (MTEB) ([Muennighoff et al., 2022](#)), evaluates 8 embedding tasks on 58 datasets covering 112 languages, providing an overview of available datasets.

This study relates most closely to [Silcock et al. \(2022\)](#), which contrastively trains an S-BERT MP-

Net model ([Reimers and Gurevych, 2019](#); [Song et al., 2020](#)) to map historical newswire articles from the same underlying article source to similar embeddings. We initialize News Déjà Vu with their model weights.

More generally, this study follows from the literature on open domain retrieval ([Karpukhin et al., 2020](#); [Thakur et al., 2021](#); [Wu et al., 2019](#)). We also draw inspiration from a large literature showing the importance of contrastive training for semantic similarity applications, which we apply to train News Déjà Vu. The anisotropic shape of the embedding space in pre-trained transformer models like BERT creates challenges for utilizing their latent features ([Ethayarajh, 2019](#)). In these models, less common words are dispersed towards the edges of the hypersphere, the sparsity of low frequency words violates convexity, and the distance between embeddings is correlated with lexical similarity. This leads to misalignment among texts with similar meanings and diminishes the effectiveness of averaging token embeddings to represent longer texts ([Reimers and Gurevych, 2019](#)). By applying contrastive training, anisotropy is mitigated ([Wang and Isola, 2020](#)), enhancing the quality of pooled sentence (or document) representations ([Reimers and Gurevych, 2019](#)).

## 3 Model Architecture and Training

The News Déjà Vu model architecture at inference time is shown in Figure 1. News Déjà Vu first recognizes and masks spans of text containing named entities (people, locations, organizations, and other miscellaneous proper nouns), as our aim is to draw parallels between articles that describe different entities in different time periods. We then replace all detected entities by the [MASK] token. Query articles are used to retrieve their semantic nearest neighbor(s) in a corpus of interest, using

<sup>1</sup><https://huggingface.co/spaces/dell-research-harvard/newsdejavu>

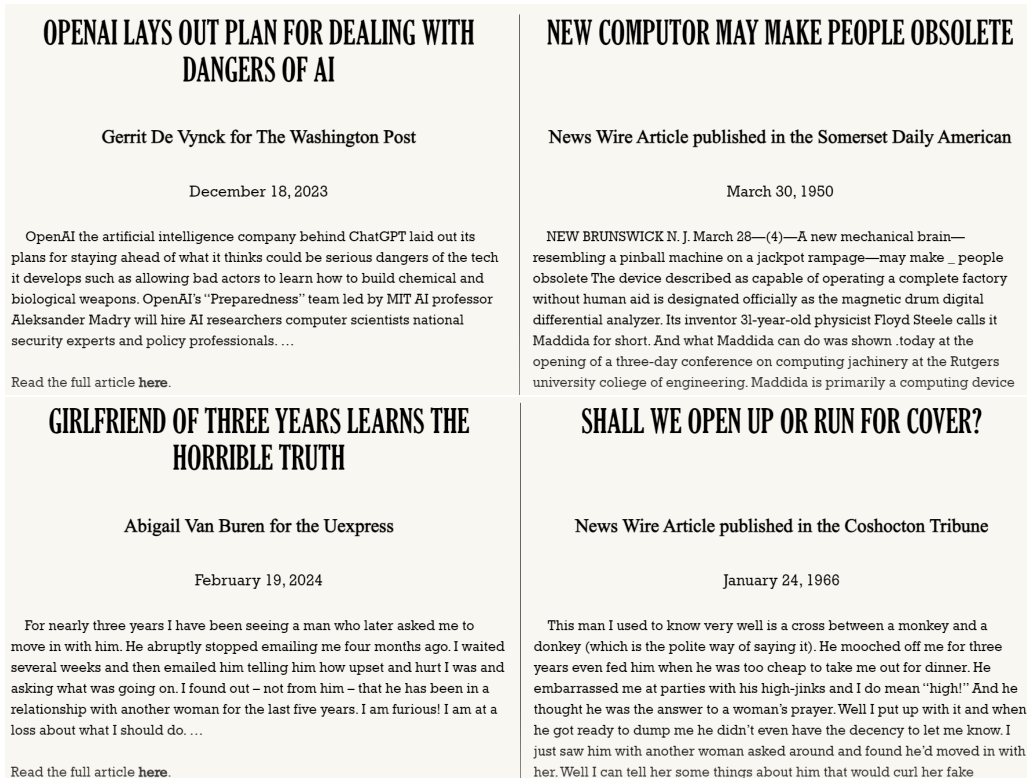


Figure 2: Examples of News Déjà Vu retrieval. The left-hand side shows a modern news article and the right-hand side shows a retrieved historical article.

the News Déjà Vu contrastively trained bi-encoder. We use the IndexFlatIP index from FAISS (Johnson et al., 2019) to perform an exact K-nearest neighbor search. Our embedding vectors are L2 normalized which makes the Inner Product metric used in the index equivalent to Cosine Similarity.

Custom training was necessary to achieve accurate NER performance with robustness to OCR noise. Table 1 describes the training data, which were drawn from randomly selected articles from off-copyright newspaper articles between 1922 and 1977. All data were double-labeled by two highly skilled undergraduate research assistants, and all discrepancies were resolved by hand. The supplemental materials contain the annotator instructions. We use the training set to fine-tune a Roberta-Large model (Liu et al., 2019). We optimised hyperparameters using Hyperband (Li et al., 2018). The best model was trained at a learning rate of  $4.7e-05$  with a batch size of 128 for 184 epochs.

Table 2 evaluates NER model performance. This model achieves an F1 of 90.4 in correctly identifying spans of text containing named entities without regards to the class, the relevant task since News Déjà Vu replaces all entities with the [MASK] token. This outperforms a Roberta-Large model

Split	Person	Org	Loc	Misc	Articles
Train	1345	450	1191	1037	227
Val	231	59	192	149	48
Test	261	83	199	181	48

Table 1: The first four columns provide the number of entities of different types in the training, validation, and test sets. The final column provides the total number of labeled articles.

finetuned on CoNLL03 by a large margin.

Model	Precision	Recall	F1
Custom NER	87.9	93.1	90.4
Roberta-Large finetuned on CoNLL03 (Conneau et al., 2019)	80.3	75.5	77.8

Table 2: Evaluation of NER models.

We would like to use News Déjà Vu for unsupervised data exploration, to retrieve historical texts that social scientists and the general public will find thought-provoking. In order to do this, we need an LLM that maps semantically similar articles to similar representations, and we found that off-the-shelf contrastively trained models - such as those described in the MTEB benchmark - did not perform satisfactorily.

Creating paired training data of modern and historical articles that do or do not have parallels would be challenging and costly. Rather, we begin with the model from [Silcock et al. \(2022\)](#), which was contrastively trained on paired historical newswire articles, with the purpose of detecting noisy duplication, rather than semantic similarity. This is a useful starting point since it has already been exposed extensively to the idiosyncrasies of historical news texts, such as OCR errors and obsolete spellings.

We further train on modern data that pairs news articles belonging to the same news story. These are drawn from Allsides, a news aggregator that collates the beginnings of articles on the same story from multiple news sites. Pairs of (the beginnings of) articles from these groupings, which typically consist of two or three texts, form positives. To create negative pairs, we used a larger pool of articles from Allsides, leveraging their pages of articles that are on the same topic, which are broader groupings than those on the same story. We embed this pool using the model from [Silcock et al. \(2022\)](#), which is a finetuned S-BERT MPNET model ([Reimers and Gurevych, 2019](#); [Song et al., 2020](#)). Then for each article that appears in a positive pair, we find the closest article (highest cosine similarity) in the pool that a) is from the same news source and b) does not appear on the same topic or story page.<sup>2</sup> Training data statistics are given in Table 3.

Training		Validation		Test		Total	
Pos.	Neg.	Pos.	Neg.	Pos.	Neg.	Pos.	Neg.
12868	12913	2757	2766	2757	2766	18,382	18,445

Table 3: Training, validation, and test sizes for the paired data used to train the retriever.

We used Hyperband ([Li et al., 2018](#)) to find optimal hyperparameters, which led us to train for 9 epochs, with a batch size of 32 and a warm-up rate of 0.39. We use S-BERT’s online contrastive loss ([Hadsell et al., 2006](#)) implementation with a margin of 0.5. The model achieves a pairwise F1 of 92.4 on the test set, outperforming models that are not finetuned, such as SBERT MPNET ([Reimers and Gurevych, 2019](#)) and [Silcock et al. \(2022\)](#), as shown in table 4.

We do not have paired modern-historical article texts for evaluation, and it is not clear how one

<sup>2</sup>In some cases there was no article that met these conditions. In these cases, we took an article from another news source.

Model	Precision	Recall	F1
Custom biencoder	93.7	91.1	92.4
<a href="#">Reimers and Gurevych (2019)</a>	83.8	85.8	84.8
<a href="#">Silcock et al. (2022)</a>	60.7	69.5	64.8

Table 4: Evaluation of biencoder models.

would create such data given the nature of the task. Rather, a skilled annotator gave 350 randomly selected modern articles and their nearest historical News Déjà Vu neighbor a short description of their main topic (typically 2-3 words). In 85.7% of cases, the modern query and its nearest historical neighbor had the same major topic. Even when the major topic is not the same, the pairs often showed other interesting parallels.

Figure 2 provides representative examples of News Déjà Vu retrieval. Other examples can be seen at [newsdejavu.github.io](#). Modern query articles are truncated due to copyright protection. Historical articles are drawn from off-copyright newspapers and reproduced in their entirety. Except for a few special editions about topics requested by our readers, these articles were selected at random, painting a representative picture of News Déjà Vu.

We have a demo where users can use their own texts to query a subset of American Stories ([Dell et al., 2023](#)), a massive scale Hugging Face dataset consisting of over 430 million historical newspaper texts.<sup>3</sup> We also make embeddings for the American Stories collection available on Hugging Face.<sup>4</sup>

## 4 The News Déjà Vu Package

The News Déjà Vu package is available on PyPI for easy install.

```
1 pip install newsdejavu
```

The package consists of the following core functionalities: `download`, `ner`, `mask`, `embed`, and `find_nearest_neighbours`. The package focuses on inference, which we expect is how the vast majority of users would like to use News Déjà Vu. For users who wish to fine-tune their own News Déjà Vu model, we recommend using LinkTransformer ([Arora and Dell, 2023](#)) or the Sentence BERT repository ([Reimers and Gurevych, 2019](#)), initializing with our pre-trained weights which are available on Hugging Face.

<sup>3</sup><https://huggingface.co/spaces/dell-research-harvard/newsdejavu>

<sup>4</sup>[https://huggingface.co/datasets/dell-research-harvard/americanstories\\_masked\\_embeddings](https://huggingface.co/datasets/dell-research-harvard/americanstories_masked_embeddings)

The download function downloads the dataset that users would like to work with. We have integrated support for American Stories (Dell et al., 2023). This command allows users to specify which states and year range they would like to download, or they can download the (very large) dataset in full. This step of the pipeline can also be skipped if the user already has a dataset that they would like to use.

```
1 import newsdejavu as ndjv
2 corpus = ndjv.download('american_stories
:1900:Alabama')
```

The ner command runs NER over the corpus.

```
1 ner_outputs = ndjv.ner(corpus, '
historical_newspaper_ner')
```

Next, these detected entities can be masked with a simple mask command and the texts embedded using News Déjà Vu with the embed command. In addition to using the News Déjà Vu model as the default, this command can also support using a local model path or downloading one from Hugging Face, for users who would like to use their own retrieval models in conjunction with the package.

```
1 masked_corpus = ndjv.mask(ner_outputs)
2 corpus_embeddings = ndjv.embed(
masked_corpus, 'same-story')
```

Users can similarly mask and embed their query article. Finally, find\_nearest\_neighbours retrieves the  $k$  closest corpus articles to the query.

```
1 dist_list, nn_list =
find_nearest_neighbours(
query_embeddings, corpus_embeddings,
k=1)
```

Users would typically like to use all these commands in sequence. The mask\_and\_embed command combines NER, masking, and embedding, and the search\_nearest\_story command combines NER, masking, embedding, and retrieval.

```
1 corpus_embeddings = ndjv.mask_and_embed(
ner_outputs)
2 nearest_articles = ndjv.
search_nearest_story(query_articles,
'historical_newspaper_ner', 'same-
story', corpus_embed =
corpus_embeddings)
```

We recommend that those who lack extensive familiarity with deep learning frameworks install it on a cloud compute service optimized for deep learning, such as Google Colab, in order to avoid the need to resolve dependencies. Tutorials on how to use News Déjà Vu on Colab will be provided on the News Déjà Vu Github Repository ([github.com/](https://github.com/dell-research-harvard/newsdejavu/)

[dell-research-harvard/newsdejavu/](https://github.com/dell-research-harvard/newsdejavu/)) and on the website.

By making semantic search an accessible tool for social scientists to apply to historical document collections, we hope to make it easier for researchers to find content that contextualizes our understanding of the parallels between past and present.

## Ethics Statement

News Déjà Vu is ethically sound. We do emphasize that it retrieves articles that use similar language, which may or may not reflect similarities in the underlying events or phenomena being described. Trained human judgement is required to draw deeper parallels between the past and present, and we hope News Déjà Vu will be a useful tool for directing researchers and the public to content of interest.

## Acknowledgements

We are grateful for research assistance from Dennis Du, Jude Ha, Alice Liu, Shiloh Liu, Stephanie Lin, Andrew Lu, Prabhav Kamojjhala, and Ryan Xia.

## References

- Abhishek Arora and Melissa Dell. 2023. Linktransformer: A unified package for record linkage with transformer language models. *arXiv preprint arXiv:2309.00789*.
- Xi Chen, Ali Zeynali, Chico Q Camargo, Fabian Flöck, Devin Gaffney, Przemyslaw A Grabowicz, Scott A Hale, David Jurgens, and Mattia Samory. 2022. Semeval-2022 task 8: Multilingual news article similarity. *ACL Anthology*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Melissa Dell, Jacob Carlson, Tom Bryan, Emily Silcock, Abhishek Arora, Zejiang Shen, Luca D’Amico-Wong, Quan Le, Pablo Querubin, and Leander Heldring. 2023. American stories: A large-scale structured text dataset of historical us newspapers. *Advances in Neural Information and Processing Systems, Datasets and Benchmarks*.
- Kawin Ethayarajh. 2019. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. *arXiv preprint arXiv:1909.00512*.

- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Ros-tamizadeh, and Ameet Talwalkar. 2018. [Hyperband: A novel bandit-based approach to hyperparameter optimization](#). *Preprint*, arXiv:1603.06560.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining ap-proach. *arXiv preprint arXiv:1907.11692*.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Emily Silcock, Luca D’Amico-Wong, Jinglin Yang, and Melissa Dell. 2022. Noise-robust de-duplication at scale. Technical report, National Bureau of Eco-nomic Research.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. *Advances in Neural Information Processing Systems*, 33:16857–16867.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Ab-hishek Srivastava, and Iryna Gurevych. 2021. Beir: A heterogenous benchmark for zero-shot evalua-tion of information retrieval models. *arXiv preprint arXiv:2104.08663*.
- Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. *International Conference on Machine Learning*, 119:9929–9939.
- Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2019. Scalable zero-shot entity linking with dense entity retrieval. *arXiv preprint arXiv:1911.03814*.

## A NER annotator instructions and results

During the NER annotation process, careful rules were developed to ensure congruence between use of labels. This appendix details those rules

### A.1 General rules

- Label the biggest span than constitutes one entity, with the exception of locations. For example “Martin Luther King High School” would be one entity, not “Martin Luther King” or “High School”.
- Label “Vietnamese government”, not just “Vietnamese”.
- “Adam Smith” is one person and the “Catholic Church” is one organization, but “Topeka, Kansas is two locations.
- If the sub-entity is ambiguous without connection to the parent entity, we label it as one entity. For instance, The Department of Electrical Engineering at Tri-State College is all one entity
- Do not include extra punctuation or spaces in the labels, unless they occur within a named entity. If OCR errors produce “First (Baptist Church”, the extra parenthesis should be included since it’s inside a named entity. But if you have “(First Baptist Church”, do not include the parenthesis.
- On a similar note, given something like “Albert Sealy’s cat”, only label “Albert Sealy”. Do not include the apostrophe and the “s”.
- Where relevant, don’t label the “the” (ie. should be “State Department” not “the State Department”).
- For newspapers and other publications/organizations use Google to see if “The” is part of their official name. For instance, “The New York Times” is the actual full name of the newspaper, so label the “The”. However, “The Bell Syndicate” doesn’t actually have “The” as part of their official name, so just label “Bell Syndicate”.
- If there is an organization/location ambiguity, we default to labeling it as a location. If the ‘location’ does an action, it is labeled as an organization. e.g. White House = LOC unless “the White House” does something. e.g. I’m going to the Natural History Museum and the Marriot hotel are both locations.
- EXCEPTION: All administrative units (countries, states) are locations and never organizations, even when doing an action.
- If a location is part of the name of an organization, we label the whole thing as an organization (e.g. “the church of christ in America” or the “Ohio agriculture department”).
- We’re defining locations as 4D: distinctive time periods that are capitalized for places are included in locations e.g. Victorian England, Ancient Greece, Nazi Germany, and Red/Communist China are all labeled as single locations.
- Cardinal direction plus location: if defined in popular speech as a specific location, label the direction also e.g. Western Europe, Central London.
- If the secondary word is plural, it is not a specific location and is not included in the label (e.g. Marriott hotels).
- When established locations are used to reference a location in relation to them, the established locations are labeled separately. E.g. “between Smith and Adams street” are two separate locations.
- This also applies to railroad lines e.g. “the Chicago-New York line” should have Chicago and New York labeled separately.
- Named objects in space (e.g. Sun, Moon, Jupiter) are miscellaneous unless someone/something is going to it or if there is a reference to a place on the object, then they are locations.

### A.2 Locations

- If a location title is used as an adjective, it is miscellaneous. e.g. “US senator” should have “US” labelled as MISC, not LOC.

### A.3 People

- For person entities, do not label Mr, Mrs, and other prefixes (such as Dr, Rev, etc.) and don’t include suffixes (e.g. “Jr.”, “III”).

- Similarly, do not label descriptors like “Deputy Sheriff” that often come before people’s names.
- Also, positions (without a name) are not named entities e.g.: The Minister of Foreign Affairs, Maharajah, the Queen should not be labelled.
- Include any nicknames when they are put in the middle of names (e.g. something like “Dwayne ‘The Rock’ Johnson” should be all labeled person).
- God is a person, but pronouns like “He” and “Him” are not.
- Animal names are miscellaneous, not person.

#### A.4 Organizations

- If an organization is used as an adjective, it is miscellaneous.
- Groups of people that go by a name (e.g.: Republicans, Cavaliers) can be orgs, but only if we’re referring to the whole group (e.g. “the republicans” referring to the whole party, not a group of five republicans).
- Liberals and conservatives are not organizations (at least in the US), since the definitions often change and are not really groups.
- Many politician names will be followed by (D/R-State). Label the D/R as miscellaneous, since the group is being used as an adjective to describe the politician. Label the state as location.
- If there is an organization/location ambiguity, we default to labeling it as a location. If the “location” does an action, it is labeled as an organization e.g. ‘the “White House” reported,’ is an organization, but ‘I’m going to the Natural History Museum’ is a location.
- All administrative units (countries, states) are locations and never organizations, even when doing an action.
- Brands of products are organizations if they are the name of the company/producer (e.g. Apple computers) or unless they are now ubiquitous (e.g. ziploc bag, polo shirt). If the brand is not the name of the company, it’s

miscellaneous. E.g.: Toyota is an org, highlander is misc, x-ray is neither, 45-caliper gun is neither.

- Ambiguous organizations are still labeled when they are non-ambiguous in context (e.g. an article talking about Cleveland saying the “rotary club”).
- This also applies to slightly ambiguous government entities like the “army”, “navy”, or “state department” - if it’s clear it’s the US entity in the article, label it as such.

#### A.5 Miscellaneous

- Adjectives derived from named entities are miscellaneous named entities e.g. nationalities used as adjectives (e.g. U.S., French, London Newspapers).
- EXCEPTION: People used as “adjectives” (e.g. “the Kennedy household”) are to be labeled as PER, not MISC.
- EXCEPTION: When a person’s name has become part of a famous location (e.g. “Eiffel Tower”, “Chandler Building”) that has its own Wikipedia page (or equivalent) the entire location is considered an entity and labeled appropriately.
- When an entity is used as a possessive, that is with an apostrophe or “of” (e.g. “Wisconsin’s cows”, “people of France”), the entity should be labeled with its original label, not MISC.
- Congressional, senatorial, constitutional are all not considered adjectives from named entities.
- Political ideologies are miscellaneous (e.g. communist, socialist, conservative, authoritarian etc.).
- Include prefixes and suffixes to these in misc labels e.g. anti-Japanese, pro-communist.
- Names of groups of people/religions are miscellaneous (only if they’re capitalized or ‘should’ be capitalized, e.g. “communists” should be misc, but “visitors” should not be).
- Titles/positions are not miscellaneous named entities e.g.: The Minister of Foreign Affairs, Maharajah, the Queen are not misc.



- Officially named initiatives/programs are misc e.g. Manhattan Project, U.S. Census.
- Names of capitalized documents or forms are miscellaneous e.g. Individual Census Reports, the Constitution.
- Distinct political acts (e.g. “Agricultural Act of 1970”) are also misc.
- Capitalized/specific names of objects are misc e.g. USS Canopus.
- Names of animals are miscellaneous e.g. Laika (the space dog).
- Events are miscellaneous.
- Events must be famous or distinct (e.g. “Pearl Harbor”, “1969 World Series”).
- Less famous events must refer to something occurring within a specific timeframe (e.g. there will be a “city council meeting” is not a named entity). The event should be unambiguous (e.g. the “city council meeting on 11/5” is not a named entity because it doesn’t specify which city council but the “Boston City Council meeting on 11/5” is).
- Christmas and other major holidays are miscellaneous.
- Brands of products are organizations if they are the name of the company/producer (e.g. Apple computers) or unless they are now ubiquitous (e.g. ziploc bag, polo shirt). If the brand is not the name of the company, it’s miscellaneous. e.g. Toyota is an org, highlander is misc but not org, x-ray is neither, 45-caliper gun is neither.
- Named objects in space (e.g. Sun, Moon, Jupiter) are miscellaneous unless someone/something is going to it or if there is a reference to a place on the object, then they are locations.

## A.6 NER Results

The following shows the breakdown of shares of different types of entities (in terms of overall tokens) from applying the NER model to newswire articles from the 20th century. The figure shows broad stability in entity mentions across time, with the exception of World War II, when location and miscellaneous entities (e.g., such as named aircraft) spike.

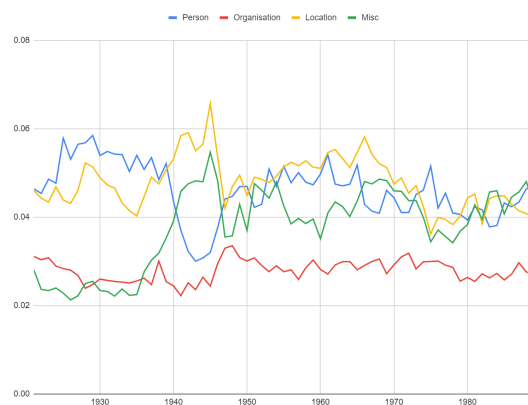


Figure 3: Shares of entity types in newswire articles.

## **B Explanation and Examples for Evaluation of News Déjà Vu**

### **B.1 Explanation**

The News Déjà Vu package was used to obtain a total of 5 historical articles from our corpus per modern article of which there were 70. These 70 articles were pulled from the websites of popular news outlets such as the Associated Press, Fox News, USA Today, and the Washington Post. This was done using the News Déjà Vu python package. Each modern article was then placed in a pair with its 5 associated historical articles. A skilled annotator was then tasked with manually classifying the articles in a modern-historical pair as being on or off the same topic.

Being on or about the same topic was defined in the following way: two articles are about the same topic if absent named entities (e.g., proper nouns), the relationship between the most important remaining concepts are essentially the same for both articles. Any two articles that fail to meet these criteria according to the skilled annotator were deemed off or not about the same topic.

Two articles being about the same topic can also be thought about in relation to two articles being on the same story and or on the same events. Topics, stories, and events are distinguished along two dimensions: time and named entities. For two articles to be about the same event, what they each describe must have occurred during the same time period (no more than a day) and have the same or closely related named entities. To be about the same story, two articles must have the same or closely related named entities, but the actions off named-entities or the relationships between named-entities may have occurred over long stretches of time. In both events and stories, named-entities serve the same function as a fictional story's cast of characters. Named-entities and when their actions occur are irrelevant for determining two articles' shared topic. Events, stories, and topics are also nested concepts. Every event is part of a story and every story is part of a topic. An article may also belong to more than one event, story, and/or topic.

An illustrative example might be "Watergate". The "Watergate Break-In" was an event that occurred on Saturday, June 17, 1972 at the Watergate Hotel in Washington, D.C. and was foiled by a security guard named Frank Willis. This is an event because it occurred on a discrete date in time with a clear set of named-entities, like the Watergate Ho-

tel, Washington, D.C., and Frank Willis. If another article was about perpetrator Virgilio Gonzalez's arrest at the Watergate Hotel that night, that article would also be about the "Watergate Break-In" because the arrest occurred at the same date in time. Virgilio Gonzalez is related to the other named entities. The Watergate Hotel in D.C. was the location of his arrest, and it was Frank Willis' tip to the police that led to his arrest. The "Watergate Scandal" would be a story as it unfolded over years and included named entities that were related to one another but did not all participate in the same events. For example, Richard Nixon did not personally break into the Watergate Hotel, but it is believed that he tried to cover up his connections to the break-in's perpetrators, like Virgilio Gonzalez. The cover-up also occurred sometime after the break-in. An article about Richard Nixon's cover-up activities would be about the same story as an article about Virgilio Gonzalez's arrest. "Political Scandals" would be an example of a topic that the "Watergate Scandal" falls under. An article about the "Watergate Scandal" and an article about President Donald Trump pressuring election officials for more votes would both be political scandals. The named entities and dates in time are different, but the most important remaining concepts, like accusations of election interference against presidents, are essentially the same.

### **B.2 Examples of Historical-Modern Article Pair Evaluation**

The following articles make five of the labeled article pairs, each containing the same modern article and one of the five historical articles retrieved using the News Déjà Vu package. Three historical articles were classified as being on the same topic and two were classified as being on different topics. In the same way that you can provide names for events, stories, or topics, like "Political Scandals", the skilled annotator was tasked with doing so for the evaluation of the modern-historical article pairs and those names are included here. In order to abide by copyright restrictions, the full modern article is not reproduced here, just a truncated version. A link to the full article is provided however. Due to OCR errors, the historical article text may appear less coherent than the modern article but still readable nonetheless.

## B.2.1 Modern Article

### Modern Article URL

<https://www.usatoday.com/story/money/food/2024/03/13/ben-jerrys-free-cone-day-2024/72944410007/>

### Modern Article Headline

Ben & Jerry's annual Free Cone Day returns in 2024: Here's when it is and what to know

### Modern Article Body

Ben & Jerry's is bringing back its annual Free Cone Day celebration this spring and is asking fans to help them beat a lofty goal.

The company wants this year's Free Cone Day to be the "biggest and best yet with 1 million scoops served," it announced Wednesday.

This year's celebration will take place on Tuesday, April 16, the company said in a news release. Free Cone Day made its return last year after a four-year hiatus.

## B.2.2 Historical Article 1

### Historical Article ID

10\_304477967-ottumwa-daily-courier-Jan-04-1943-p-1.jpg

### Historical Article Headline

Additional Cut In Ice Cream Output

### Historical Article Body

Washington, D. C.-(P)-The war production board today limited January production on ice cream to 50 per cent of the amount each manufacturer made last October.

- This represented a reduction of one sixth from December, when each manufacturer was permitted to make 60 per cent of his October amount.

The order also applies to frozen custard, milk sherbet, other frozen desserts and ice cream mix.

W.P.B. said the order "was issued at the request of United States department of agriculture . . . » to further relieve the butter shortage."

## On Topic with Modern

True

### Topic Name / Notes

Ice Cream

## B.2.3 Historical Article 2

### Historical Article ID

1\_106022053-titusville-herald-Apr-11-1939-p-1.jpg

### Historical Article Headline

52,299 Jam White House Lawn 'For Annual Easter Ego Rolling

### Historical Article Body

OY 2£2RE NAAGCIAER Pear

WASHINGTON, April 10.—Even Alice in Wonderland couldn't have dreamed anything like the Easter rolling which took 52,259 children and adults—by the gate keepers' count—to the private grounds of the White House today.

Through the gates streamed a ten-year-old boy dressed like a white rabbit with flapping pink ears... a live white rabbit scampering on a leash . . . & huge chocolate bunny perched beside a grinning infant in a baby carriage.

Bands played. A magician did tricks. Crowds gathered about a Punch and Judy show. Everywhere, - there were children and eggs—eggs crunched underfoot, smeared on young faces, salling between stately old White House trees and rolling down slopes.

It was like a glorified country picnic. Thousands spread their lunches on the ground. Then it wasn't, for there on the White House porriica wag Lhe President of the United States, waving and wishing, he said, that he could "be down there with you."

And four times during the pleasant sunny day the country's gracious First Lady appeared. Three times she made trips through the grounds, smiling, waving and calling "how-do-you-do" to those along her path. | As she walked she saw children crawl; about on all fours, rolling themselves down slopes, sitting spraddie-legged on the ground to juggle eggs and swingIng from iow limbs of trees.

It was the seventh time, the President pointed out, that he and Mrs. Roosevelt had entertained the children at the White House lawn. The roll has been an annual event at the White House, except for the World war years, since 1878.

"It's a wonderful day," he declared before he disappeared into the house, "and I hope you u enjoy yourselves very, much,?"

### On Topic with Modern

False

### Topic Name / Notes

one is about ice cream and the other is about Easter

### B.2.4 Historical Article 3

#### Historical Article ID

1\_10883667-daily-globe-Mar-01-1949-p-1.jpg

#### Historical Article Headline

Butter vs. Oleo Battle, Annual Congress Feature, Starts Today

#### Historical Article Body

Washington — çYfj—The butter vs. oleo battle, an annual feature of congress, got underway today,

Rep. Poage (D-Tex) fired the first shot with an attack on "the butter lobby." He said it is more interested in tryilng "to hill yelJow margarine" than it is in protecting consumers from fraud.

Poage was the first witness at house agriculture committee bhearmgs on 29 bills regarding margarine, Some of them would repeal the taxes on margarine, but prohibit the manufacture und sale of yellow margarine,

Poage sald this latter proposal is an effort en the part of the butter industry "to appear to yield to outraged public sentiment against inexcusable favoritism In favor of one wholesome food

-aguinst anolher without actually lgiving up any af the special pri vilige butler" bas so long enjoy: ed."

#### POAGE'S BILL

Poage has a bill which would lift, federal taxes on margarine' and permit manufacturers 10 col: or it yellow. But it would require public eating places to notify cus: tomers \$f they serve margarine.

The issue as whether there should be a federal tax on the

l butter substitute. The government ;Now taxes all of it, with on ex. tra toll if at is ecalored

President Truman for repeal of the margarine ti ;@s, and so did the Democri platform on which he ran list fall, The Republican platform did ; not mention the subject.

Tf the controversy follows its new almost traditional form — butter partisans will suggest that margarine could be colored anything from purple or green to a bright, cherry red=anything but yellow,

The margarine folks will say butter has no exclusive cham to yellow, and in fact the yellow color has to be added ta some bulter.

Southern Democrats from rural districts and northern Democrats sand even Republicans — from big city districts will ally for the moment in support of marparine,

There are two types of Dills before ihe commiltee,

The most numerous would abolish all federal taxes on margnrine. Others would but out st least the added nx on colored margare.

At the other end of the line are bis by Reps. August 11. An ;aresen (R-Minn), und Granger '(D-Uiah) to prohibit the manufacture and sule of yellow marigaring, STARTED 69 YEARS AGO

The whole thing started 63 years ago, when the first federal anti-margarine Jaw was pussued. That law and others which followed have been under attack almost every year since.

The oleo forees made their greatest progress last yeat, when a bill fo repeal the margarine taxes passed the house handily. It never got to a vote in the sen l ate.

A federal tax of 10 cents a pound is now paid on all colored margarine sold at retail. The retail tax on the uncolored product is one-quarter cent a pound.

Margarine manufacturers pay a federal tax of \$600 a year; wholesalers pay a tax of \$480 4 year Sf they handle colored margarine, and \$200 a yenr for uncolored, and retailers are taxed \$48 n year fo handle colored mare garine and \$6 10 handle uncol

### On Topic with Modern

True

### Topic Name / Notes

food

## B.2.5 Historical Article 4

### Historical Article ID

26\_90193365-morning-herald-Apr-04-1946-p-1.jpg

### Historical Article Headline

No More KP In The Air Forces Of Army

### Historical Article Body

Washington, April 3 (P) There will be no more KP (Kitchen Police) duty in the Army Air Forces under a new program announced today.

Soldiers will still peel spuds and wash dishes, But those who do will be permanently assigned to the task and will be called "mess attendants." The announcement adds that they "will be afforded an opportunity to make an Army career of food service,"

The old system of assigning all men on the roster to KP in turn is being abolished.

The AAF announcement said that "many (local exigencies and personnel problems" prevent setting a definite date for the establishment of Utopia.

### On Topic with Modern

True

### Topic Name / Notes

Food

## B.2.6 Historical Article 5

### Historical Article ID

5\_233435376-circleville-herald-Dec-30-1976-p-1.jpg

### Historical Article Headline

International Falls No Sunny Spa

### Historical Article Body

INTERNATIONAL FALLS, Minn. (AP) — At the close of a year, a time for reflection, hardly a place is better suited than this for that worthy exercise.

There is something to be learned from this place, something other than what everybody already knows from the nightly weather report: that it is the coldest place in the 48 contiguous states.

When winter's fangs bite into this little spot on the Canadian border — in the first half of this month the thermometer managed to hang above zero for only four brief hours — living becomes an adventure and humility a daily lesson. Nature's elemental severity invites long thoughts about man's standing in the Great Scheme.

"I think we worry more about the simple necessities of survival than most people do," said Frank Bohman, a philosophic aviator who has lived here all his 52 years.

"Having enough food in the house, enough fuel, a backup heating system, these are real concerns. I would imagine that in gentler climates people take survival for granted."

For the record, when the earth tilts toward winter, winds borne on the jet stream sweep from the North Pole down the interior flank of the Canadian Rockies and pivot eastward right at this point, so that the average yearly temperature here is 37.5 degrees and the annual snowfall 50 inches. Readings in the minus 30s and 40s are commonplace during the winter.

The cold grips so fiercely, in fact, that it all but refuses to let go. The ground freezes five feet down, untillable until June.

The town is on the granite shore of Rainy Lake, one of creation's masterpieces, a 340-square-mile work of art done in a freeform of coves and bays and flecked with 1,600 tiny granite islands timbered with pine.

Thus in the summertime the area is awash with tourists, regulars who return to their summer places on the islands, weekenders seeking walleyed pike and clean air, visitors with cameras, water skis and time to make the two-hour drive up from Duluth.

When the summer crop of frolickers is harvested, however, only a bold band of the hearty remain to face the long dark winter.

That yearly experience has given them a palpable sense of neighborliness, a closeness such as a shared secret brings. When the ice breaks up each May they have earned a communal handshake that says nice going everybody, we did it again, we didn't quit.

Those brave souls number 9,109 in. International Falls and the nearby communities of South International Falls and Ranier. About the same number

live across the Rainy River Bridge at Fort Frances, Canada.

The spirit of hands-across-the-sea, or in this case, the river, comes naturally; nature's legacy knows no international boundary. Indeed, one longtime Chamber of Commerce president in International Falls, Gordy McBride, was a Canadian citizen.

**On Topic with Modern**

False

**Topic Name / Notes**

one about ice cream and the other about weather