

# 1 **Supplementary Materials**

## 2 **1 Code Availability**

3 All code used to create Newswire is available at our public Github repository.

## 4 **2 Model Details**

### 5 **2.1 Digitalization**

6 To identify content regions like articles, headlines, and ads in newspaper scans, we employed YOLOv8  
7 (Medium) (20), starting from the official YOLOv8m pretrained model. We trained for over 100  
8 epochs on 2,202 scans with 48,874 layout objects. This model achieves a 0.91 mAP50:95 for articles  
9 and 0.84 mAP50:95 for headlines. The confidence threshold was lowered to 0.1 to enhance recall.

10 Text bounding boxes were classified as legible, borderline, or illegible using MobileNetV3 (Small)  
11 (10), initialized from a PyTorch Image Models checkpoint (22). Training involved 979 examples  
12 (678 legible, 192 borderline, 109 illegible) over 50 epochs with weighted Cross Entropy Loss and a  
13 learning rate of  $2e-3$ , which was reduced every 20 epochs, as described in (6).

14 Text prior to 1920 was OCRd using EffOCR (5; 3), while later text was transcribed using Tesseract.

15 Articles with texts spanning multiple bounding boxes need to be associated, for which we use a  
16 customized RoBERTa cross-encoder model to predict whether the content in one bounding box  
17 continues the content in another, as described in (17).

### 18 **2.2 De-duplication of content**

19 To detect reproduced content, we use a bi-encoder model, as developed by (18). This model, based  
20 on the S-BERT MPNET model (16; 19), is fine-tuned to learn similar representations for reproduced  
21 articles and dissimilar representations for non-reproduced ones. The model is fine-tuned on a manually  
22 labeled dataset of near-duplicate articles. The model is fine-tuned with a learning rate of  $2e-05$ , using  
23 the online contrastive loss implementation from S-BERT (9), for 16 epochs, with a batch size of 32  
24 and a 100% warm-up. A cosine similarity metric with a 0.2 margin is used. The model achieves an  
25 ARI of 91.5 on a hand-labeled test set of over 100 million pairs of articles (54,996 positives pairs and  
26 100,914,159 negative pairs). This performance is compared to other baselines in table ?? in the main  
27 text.

28 Once the embeddings are generated, to create clusters of near duplicates we apply highly scalable  
29 single-linkage clustering, setting a cosine similarity threshold of 0.94. Articles are represented as  
30 nodes within a graph, and edges are formed when the cosine similarity surpasses the threshold. Edge  
31 weights are calculated based on the negative exponential of the time gap (in days) between the articles.  
32 To ensure clusters are meaningful, we use Leiden community detection, which helps mitigate the risk  
33 of false-positive edges merging unrelated articles into the same cluster.

34 To further refine the clustering, we exclude any clusters containing more than 50 articles that span  
35 over five different dates. Similarly, clusters with more than 50 articles are removed if the number  
36 of articles exceeds twice the count of unique newspapers from which they originate. These criteria  
37 ensure the exclusion of clusters that, although correctly grouped based on a shared source, are not  
38 useful for the Newswire dataset.

39 A detailed analysis of errors is provided in (18). Typically errors are articles about the same story  
40 from different wire sources, or updates to a story as new events unfolded.

### 41 **2.3 Detection of wire content**

42 To accurately filter out non-wire content, we fine-tuned a Distil-RoBERTa classifier on a hand-labeled  
43 training set of 1,459 samples. The model was trained for 20 epochs with a batch size of 64 and a

44 learning rate of  $5e-5$  with an AdamW optimizer. All hyperparameters were selected based on the  
45 model’s performance on a validation set containing 336 labeled samples. The final model achieved  
46 an F1 of 0.96 on a test set containing 448 samples.

## 47 **2.4 Georeferencing**

48 Our georeferencing pipeline consists of multiple steps designed to extract the dateline from each  
49 cluster of reproduced articles. As a first step, we train a DistilBERT classifier to detect bylines from  
50 each article on a training set of 1,392 hand-labeled samples. The model was trained for 25 epochs  
51 with a batch size of 16 and a learning rate of  $2e-5$  with an AdamW optimizer. All hyperparameters  
52 were selected based on the model’s performance on a validation set containing 464 labeled samples.  
53 The final byline classifier achieved an F1 of 0.92 on a test set containing 464 samples.

54 For each article within a given cluster, we take all possible  $n$ -grams from the detected bylines,  
55 matching each consecutive sequence of words to GeoNames’ dictionary of city and country names.  
56 We additionally detect state names and state abbreviations within bylines. We first search for matches  
57 among capitalized  $n$ -grams, as most datelines in our corpus are capitalized, searching across all  
58  $n$ -grams only in the event that we do not find a match.

59 Once we have potential matches for each article in a cluster, we aggregate these matches to get a  
60 tentative match for the city, state (if one exists), and country in each cluster dateline. For both state  
61 and country, we take the most common potential match across all articles in the cluster. As some city  
62 names may be substrings of other city names (for example, York and New York), we additionally  
63 weight the count of each potential city match by a function of the length of the city name. In all cases,  
64 if the tentative match fails to appear in at least 15% of all articles in the cluster, we proceed without a  
65 tentative match; this is to prevent the pipeline from detecting errant place names in clusters with no  
66 dateline. The AP stylebook additionally designates a list of 56 cities which are allowed to appear in  
67 AP articles without an associated state/country name – to address these cases, we manually match  
68 these cities to their associated states/countries.

69 Having a tentative match for the city, state, and country in which each article cluster was written, we  
70 attempt to merge these tentative matches with GeoNames’ dataset of all cities with a population of at  
71 least 500 residents. Some datelines that contain locations other than cities, such as the Johnson Space  
72 Center, or very sparsely populated areas may fail to be matched as a result of this process. After  
73 running the georeferencing pipeline over our entire sample, we manually inspected the matches for  
74 any particularly common instances of these non-city datelines. We include further explanation of  
75 these exceptions in the “wire\_location\_notes” field associated with the cluster.

76 On a test set of 2,324 hand-labeled georeferenced clusters, we find that the pipeline has an accuracy  
77 of 94.9%.

### 78 **2.4.1 Benchmarking against GPT-4o-mini**

79 We additionally benchmark our georeferencing pipeline against GPT-4o-mini, passing in the following  
80 prompt:

81 I will feed you the beginning snippet of multiple articles belonging to a given cluster – in a cluster,  
82 articles should all be the same. If there is a geographic byline belonging to the articles in a cluster, I  
83 would like you to output the location. If it is in the United States, please give me the city name, state,  
84 and country. If it is not in the United States, please give me both the city name and the country name.

85 Some articles in the cluster may have a byline while others may not – if there are multiple different  
86 locations, please output only the one you believe is correct. Only output locations that correspond to  
87 the article byline – if there are other articles mentioned in the text but that are not part of a byline,  
88 ignore these. Please output only a single location and nothing else. If there is no location, output  
89 None.

90 For example, the following snippet: “In Vienna, Austria, there is much indignation because in the  
91 Balkan states a monument has been erected in honor of the student” has no location in the byline –  
92 Vienna does not belong to a byline. You should output None.

93 Meanwhile, the snippet: “LOS ANGELES, Jan. 27.–The appeal of Alexander Pantages to his  
94 conviction on charges of having assaulted” has Los Angeles in the byline. You should output Los  
95 Angeles, California, United States.

96 Remember, please output only a single location and nothing else. If there is no location, output  
97 None.

98 The above prompt achieves an accuracy of 85.3% on the same test set of 2,324 hand-labeled  
99 georeferenced clusters, compared to our pipeline’s accuracy of 94.9%.

## 100 2.5 Topic tagging

101 Two types of topics are tagged in the dataset.

102 First, we tag topics of particular interest during this period (Politics, Crime, Labor movement,  
103 Government regulation, Protests, Civil rights, Antitrust). To create training data for these models, we  
104 developed a pipeline to efficiently extract articles, as random sampling would not lead to many on  
105 topic articles. We did this in two steps. First, we trained a BART-large (13) bi-encoder on MNLI (23),  
106 using the Dense Passage Retrieval (DPR) infrastructure (12). We trained for 40 epochs, with a batch  
107 size of 32, and a learning rate of  $7e-05$ . This is a re-ranking model, so at inference time, it ranks all  
108 embedded texts with respect to a query text. We embedded all Newswire articles with this model,  
109 and formatted queries as “this example is about topic” (e.g., “this example is about civil rights”).  
110 From the results, we extracted the highest scoring articles. We run zero-shot classification (using  
111 Huggingface’s implementation, based on bart-large-mnli (13)) to classify whether these texts were  
112 on topic or not, compared to the same query. We then sampled from the on-topic and not on-topic  
113 predictions to create our datasets. These datasets were then manually labeled. For each topic we  
114 then trained a binary topic classifier. Table 1 gives hyperparameters for each model. The size of the  
115 labeled datasets and the evaluation results on the test set are shown table in the main text.

116 The second type of classifier is a multi-class classifier, which categorises data into the classes from  
117 the Comparative Agendas project (2) (30 major policy topics, such as Labor, Immigration, and  
118 Employment, Education, Environment, Energy, Immigration, Transportation). To train this, we use  
119 data from the Comparative Agenda project, as they have already labeled 4,026 short article synopses  
120 from the New York Times according to these policy topics. As we wanted to train on articles, not  
121 on synopses, we use a semantic similarity model (S-BERT MPNet) to match these synopses to the  
122 articles that they are summarising. We are able to match 1847 articles, and this match has a top-1  
123 retrieval accuracy of 95%, evaluated over 44 articles. These 1847 articles form our training data for  
124 this multi-class classifier. We used these to fine-tune a RoBERTA-large model, for 4 epochs with a  
125 batch size of 32 and a learning rate of  $5e-5$ . We found that the results of this classifier for four topics  
126 (sports, fires, weather and natural disasters, and death notices) were poor, due to a small amount  
127 of labeled data. So in these four cases, we replaced the labels with the results of binary classifiers  
128 trained on these topics, using the same process as for the other binary classifiers. The results of this  
129 second classification process were evaluated on randomly selected hand-labeled Newswire articles,  
130 with an accuracy rate of 87%.

## 131 2.6 Named Entity Recognition

132 Off-the-shelf NER did not perform satisfactorily on this data, so we trained a custom model. For  
133 training data, we randomly selected articles from Newswire, which were hand-labelled. These data  
134 are described in table 2. All data were double-labeled by two highly-trained undergraduate research  
135 assistants, and all discrepancies were resolved by hand. Annotator instructions are reproduced in  
136 full in (7). We used these to fine-tune a Roberta-Large model (14) for 184 epochs, with a batch size

Topic	Base model	Learning rate	Batch size	Epochs
Politics	RoBERTa-large	1e-6	8	50
Crime	RoBERTa-large	1e-6	8	50
Labor movement	distilRoBERTa-base	1e-5	32	50
Government regulation	RoBERTa-large	5e-6	8	50
Protests	distilRoBERTa-base	1e-5	32	50
Civil rights	RoBERTa-large	1e-5	8	50
Antitrust	RoBERTa-large	1e-5	8	50
Sports	RoBERTa-large	1e-6	8	50
Fires	distilRoBERTa-base	5e-6	16	30
Weather and natural disasters	distilRoBERTa-base	5e-6	16	30
Death notices	RoBERTa-large	1e-5	8	50

Table 1: Topic classifier training details

137 of 128, and a learning rate of 4.7e-05. Table 2 describes the training data and performance. These  
138 results are benchmarked in table ?? in the main text.

Entity Type	Data			Evaluation		
	Train	Eval	Test	Precision	Recall	F1
Location	1191	192	199	87.4	94.5	90.8
Misc	1037	149	181	73.7	68.6	79.6
Organisation	450	59	83	80.7	80.7	80.7
Person	1345	231	261	92.9	95.8	94.3

Table 2: NER data and performance

## 139 2.7 Entity Disambiguation

140 To disambiguate entities to Wikidata/Wikipedia we start with the NER output and subset it to [PER]  
141 (person) tags since we are most interested in them. We then collect each named entity within and  
142 across all newspaper articles on a given day and run it through our customized entity coreference  
143 pipeline to collapse all entity mentions on a given day into a single prototype (cluster of mentions).  
144 We use this prototype to disambiguate the constituent mentions to the entity’s Wikidata ID.

145 We imagine entity coreference and disambiguation as semantic textual similarity tasks. Entity  
146 coreference can be seen as linking similar entity mentions, and disambiguation as linking an entity  
147 mention to a template created by Wikipedia and Wikidata. The template is constructed using  
148 the entity’s name, alias, and occupation from Wikidata and concatenating it with the entity’s first  
149 paragraph in Wikipedia. Semantic similarity is measured by information that is encoded by custom  
150 contrastively trained bi-encoder models based on Sentence Transformers (16).

151 We process a Wikipedia XML dump <sup>1</sup> from November 11, 2022, and collect mentions of each entity  
152 (that appears as a hyperlink in the dump). We then split entities into a train-test-val split and pair up  
153 mentions of the same entity and associated context (defined by the paragraph containing the entity  
154 mention). These are positive pairs. We pair up an entity mention with mentions of another entity  
155 to form ‘easy negatives’. We augment our training data by adding ‘hard’ negatives where we use a  
156 novel approach of using disambiguation pages from Wikipedia that contain confusables of popular  
157 entities in the Wikiverse. For instance, the disambiguation page "John Kennedy" contains, John F.  
158 Kennedy the president, John Kennedy (Louisiana politician) (born 1951), a United States Senator  
159 from Louisiana, and John F. Kennedy Jr. (1960–1999), son of President Kennedy. We sample some  
160 contexts where John F. Kennedy was mentioned and pair them up with a context around a mention  
161 of an entity within a disambiguation page and treat this as a hard negative pair. We found that the  
162 performance of our models improved a lot by having a decorator or a set of special tokens ( $[M]$  Entity  
163  $[\backslash M]$  around an entity mention (24). For example, consider this context about President Kennedy

<sup>1</sup><https://dumps.wikimedia.org/>

164 "Eisenhower sharing a light moment with President-elect [M] John F. Kennedy [M] during their  
165 meeting in the Oval Office at White House". Some contexts naturally have multiple entities, like  
166 "Eisenhower" and "John F. Kennedy" in this case. We found that we can improve the features of these  
167 special tokens by further augmenting our training data with in-context negatives - pairing up these  
168 contexts with multiple negatives that only differ in the placement of the special tokens. With all of  
169 the variants ready we have, 179069981, 5819525, and 5132565 train, val, and test pairs respectively.  
170 We use a sequence length of 256 and truncate contexts around the mentions when necessary. We start  
171 with an *all-mpnet-base-v2* model sourced from the Hugging Face hub (21) and fine-tune it using these  
172 pairs. We train the model in Pytorch (15) with hyperparameters tuned with hyperband implemented  
173 within Weights and Biases (1).

174 We use Online Contrastive Loss as implemented in (16) and use AdamW as the optimizer with a  
175 linear warmup scheduler (20%). We train on 4 Nvidia A6000 GPUs with a batch size of 512, a  
176 learning rate of 1e-5, and a contrastive margin of 0.4. We run it for only a single epoch - seeing each  
177 pair in the train split only once. The best model is selected using pair-wise classification F1 on the  
178 validation set (the best val F1 was 92.75%). With a large dataset like this, we found it useful to divide  
179 it into 10 chunks before we began training. After finishing each chunk (1/10 of an epoch), since we  
180 resumed training on an intermediate checkpoint, we lowered the learning rate to 2e-6 after the first  
181 chunk, to reduce the chances of the optimizer overshooting the minima. Because training each chunk  
182 started with a warmup, effectively, our strategy simulated a linear scheduler with restarts.

183 Once the model is trained we embed all the newspaper articles and cluster the embeddings of articles  
184 printed on the same date using Hierarchical Agglomerative Clustering implemented with Scikit-Learn  
185 (4) with average linkage, cosine metric, and a threshold of 0.15. The clusters from this exercise are  
186 essentially mentions of the same entity on a given day. We average the embeddings within a cluster  
187 to create entity prototypes for each date. We will use these prototypes for disambiguation.

188 Next, we prepare a lookup corpus for disambiguating entity mentions (or prototypes) to the right entity  
189 using semantic information from both the context around the mention and information about it from  
190 a template we create. To create the template, we obtained names, aliases, and occupations/positions  
191 held by individuals from Wikidata. Consider the example of President Kennedy - "John F. Kennedy  
192 is of type human. Also known as Kennedy, Jack Kennedy, President Kennedy, John Fitzgerald  
193 Kennedy, J. F. Kennedy, JFK, John Kennedy, John Fitzgerald "Jack" Kennedy, and JF Kennedy. Has  
194 worked as politician, journalist, statesperson". We then suffix this template with the first paragraph of  
195 the associated Wikipedia page.

196 Next, we adapt our coreference model for the disambiguation task. We link up the contexts with entity  
197 mentions with the associated entity template to form positive pairs. Easy negatives link contexts with  
198 random entity templates. As with our coreference training, we utilize Wikipedia disambiguation  
199 pages and family information from wiki data to associate entity contexts with hard negative templates.  
200 We then split entities in an 80-10-10 train-val-test split ending up with 4202145, 522385, and 528709  
201 pairs in the respective split. We fine-tune our coreference model with similar hyperparameters as  
202 the coreference training, except without restarts (or chunking) and with the learning rate of 2e-6,  
203 batch size of 256, and 20% warmup. The model was trained for 1 epoch and the best checkpoint was  
204 selected using classification F1 as before (max validation F1 was 97%). Since the disambiguation  
205 of newspapers to the knowledge base is our main task, we adapt the training domain further to  
206 newspapers. We prepare a gold dataset to fine-tune the model on pairs crafted from newspaper  
207 contexts and Wikipedia templates. First, we obtained the names and aliases of individuals from  
208 Wikidata. Then, we search for them in our newspaper corpus, hand labeling whether they refer to  
209 the person searched for. When they do not match, these form hard negatives. We form extra hard  
210 negatives by matching an entity with another entity mentioned in the same context. We also form  
211 Wikipedia hard negatives by matching an entity with another entity mentioned in the same Wikipedia  
212 disambiguation dictionary. Finally, we create easy negatives by matching with a random entity. This  
213 dataset is described in table 3. We start with the model trained on Wikipedia pairs and fine-tune the  
214 model with an identical training setup. The maximum validation F1 achieved was 85%.

Split	Positives	Easy negatives	Hard Negatives	Wikipedia hard negatives
Train	1426	1299	1460	861
Eval	189	175	184	118
Test	198	180	183	130

Table 3: Data for finetuning entity disambiguation

215 At inference time, we prune our knowledge base to remove extraneous entities. First, we only keep  
216 those entities that have either a birth or a death date. Second, we only keep those people born before  
217 1970 (considering the period of our data). If the birth date was missing, the entity was retained.  
218 Finally, we remove those entities having no overlap and a high edit distance between the Wikidata  
219 label and the associated Wikipedia page’s title - this allows us to keep only those Wikidata entities  
220 whose Wikipedia page corresponds to the actual entity and not something related to it. Our pruning  
221 exercise brings the total number of entities in our knowledge base from 1.8 million to about 1.12  
222 million. We then embed the templates of these entities using our fine-tuned disambiguation model  
223 and stored them in an FAISS IndexFlatIP index (11). Since our embeddings are normalized, Inner  
224 Product boils down to Cosine Similarity. We then use the date-entity clusters obtained before and  
225 embed the mentions within each cluster using the model trained for disambiguation, average them  
226 (within-cluster), create entity-date prototype embeddings, and treat them as queries. To improve  
227 the quality of our results, we utilize Qrank <sup>2</sup> which ranks Wikidata entities by aggregating page  
228 views on Wikipedia, Wikispecies, Wikibooks, Wikiquote, and other Wikimedia projects. We first  
229 retrieve the 10 nearest neighbors of each query. We keep only those neighbors that are at most 0.01  
230 Cosine Distance away from the nearest match. We then use Qrank to rerank these results, essentially  
231 preferring the popular entity in cases where the returned matches are very close to each other. The  
232 Wikidata ID of the nearest embedding (after re-ranking) is then assigned to the date-entity cluster  
233 associated with the query, essentially disambiguating the clusters as well as their constituents to  
234 Wikidata. This of course is akin to treating disambiguation as a semantic retrieval problem and not  
235 handling out-of-knowledge-base entities. Our architecture allows us to use the Cosine Similarity  
236 between the entity-date prototype and the nearest template to evaluate whether or not the entity is an  
237 acceptable match. Anything lower than the threshold can be considered as either an incorrect match  
238 or out of the knowledge base. We tune a no-match threshold using a sample of human-annotated data  
239 from the Newswire. We annotate the output of our disambiguation pipeline on a set of 6,425 pairs  
240 sampled from 13 years - as correct if the returned entity is correct and incorrect when it is not. We  
241 then find the cut-off threshold that maximizes pair-wise classification precision and use that as the  
242 no-match threshold.

<sup>2</sup><https://github.com/brawer/wikidata-qrank/tree/main>

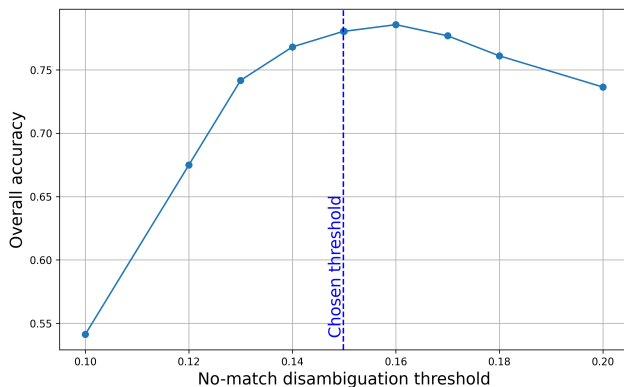


Figure 1: Sensitivity of disambiguation results to choice of no-match threshold.

## 243 2.8 Models and Dataset

244 We have made our models (see Table 4) and training/evaluation data available on the Hugging Face  
245 hub for reproducibility and ease of access by other practitioners.

Repo Name	Content
dell-research-harvard/NewsWire	The Newswire dataset
dell-research-harvard/historical_newspaper_ner	NER model for Historical Newspapers
dell-research-harvard/LinkMentions	Coreference model trained on Wikipedia
dell-research-harvard/LinkWikipedia	Disambiguation model trained on Wikipedia
dell-research-harvard/NewsLinkWikipedia	Disambiguation model fine-tuned on newspapers
dell-research-harvard/topic-politics	Topic model for politics
dell-research-harvard/topic-crime	Topic model for crime
dell-research-harvard/topic-labor-movement	Topic model for the labor movement
dell-research-harvard/topic-govt-regulation	Topic model for government regulation
dell-research-harvard/topic-protests	Topic model for protests
dell-research-harvard/topic-civil-rights	Topic model for civil rights
dell-research-harvard/topic-antitrust	Topic model for antitrust
dell-research-harvard/topic-sports	Topic model for sports
dell-research-harvard/topic-fires	Topic model for fires
dell-research-harvard/topic-weather	Topic model for weather and natural disasters
dell-research-harvard/topic-obits	Topic model for death notices
dell-research-harvard/byline-detection	Byline detection model
dell-research-harvard/wire-classifier	Classifier for wire articles

Table 4: Models and Dataset on the Hugging Face Hub

## 246 3 Dataset Information

### 247 3.1 Dataset URL

248 Newswire can be found at [https://huggingface.co/datasets/dell-research-harvard/](https://huggingface.co/datasets/dell-research-harvard/newswire)  
249 [newswire](https://huggingface.co/datasets/dell-research-harvard/newswire).

### 250 3.2 DOI

251 The DOI for this dataset is: 10.57967/hf/2423.

### 252 3.3 Metadata URL

253 Croissant metadata for Newswire can be found at [https://huggingface.co/api/datasets/](https://huggingface.co/api/datasets/dell-research-harvard/newswire/croissant)  
254 [dell-research-harvard/newswire/croissant](https://huggingface.co/api/datasets/dell-research-harvard/newswire/croissant).

### 255 3.4 License

256 Newswire has a Creative Commons CC-BY license.

### 257 3.5 Dataset usage

258 The dataset is hosted on huggingface, in json format. Each year in the dataset is divided into a distinct  
259 file (eg. 1952\_data\_clean.json).

260 An example from Newswire looks like:

```
261 {  
262     "year": 1880,  
263     "dates": ["Feb-23-1880"],
```

```

264     "article": "SENATE Washington, Feb. 23.--Bayard moved that in respect of the
265         memory of George Washington the senate adjourn ... ",
266     "byline": "",
267     "newspaper_metadata": [
268         {
269             "lccn": "sn92053943",
270             "newspaper_title": "the rock island argus",
271             "newspaper_city": "rock island",
272             "newspaper_state": " illinois "
273         },
274         ...
275     ],
276     "antitrust": 0,
277     "civil_rights": 0,
278     "crime": 0,
279     "govt_regulation": 1,
280     "labor_movement": 0,
281     "politics": 1,
282     "protests": 0,
283     "ca_topic": "Federal Government Operations",
284     "ner_words": ["SENATE", "Washington", "Feb", "23", "Bayard", "moved", "that",
285         "in", "respect", "of", "the", "memory", "of", "George", "Washington",
286         "the", "senate", "adjourn", ... ],
287     "ner_labels": ["B-ORG", "B-LOC", "O", "B-PER", "B-PER", "O", "O", "O", "O",
288         "O", "O", "O", "O", "B-PER", "I-PER", "O", "B-ORG", "O", ...],
289     "wire_city": "Washington",
290     "wire_state": "district of columbia",
291     "wire_country": "United States",
292     "wire_coordinates": [38.89511, -77.03637],
293     "wire_location_notes": "",
294     "people_mentioned": [
295         {
296             "wikidata_id": "Q23",
297             "person_name": "George Washington",
298             "person_gender": "man",
299             "person_occupation": "politician"
300         },
301         ...
302     ],
303     "cluster_size": 8
304 }

```

305 The data fields are:

- 306 - year: year of article publication.
- 307 - dates: list of dates on which this article was published, as strings in the form mmm-DD-YYYY.
- 308 - byline: article byline, if any.
- 309 - article: article text.
- 310 - newspaper\_metadata: list of newspapers that carried the article. Each newspaper is repre-
- 311 sented as a list of dictionaries, where lccn is the newspaper's Library of Congress identifier,
- 312 newspaper\_title is the name of the newspaper, and newspaper\_city and newspaper\_state
- 313 give the location of the newspaper.



314 - `antitrust`: binary variable. 1 if the article was classified as being about antitrust.

315 - `civil_rights`: binary variable. 1 if the article was classified as being about civil rights.

316 - `crime`: binary variable. 1 if the article was classified as being about crime.

317 - `govt_regulation`: binary variable. 1 if the article was classified as being about government  
318 regulation.

319 - `labor_movement`: binary variable. 1 if the article was classified as being about the labor movement.

320 - `politics`: binary variable. 1 if the article was classified as being about politics.

321 - `protests`: binary variable. 1 if the article was classified as being about protests.

322 - `ca_topic`: predicted Comparative Agendas topic of article.

323 - `wire_city`: City of wire service bureau that wrote the article.

324 - `wire_state`: State of wire service bureau that wrote the article.

325 - `wire_country`: Country of wire service bureau that wrote the article.

326 - `wire_coordinates`: Coordinates of city of wire service bureau that wrote the article.

327 - `wire_location_notes`: Contains wire dispatch location if it is not a geographic location. Can  
328 be one of “Pacific Ocean (WWII)”, “Supreme Headquarters Allied Expeditionary Force (WWII)”,  
329 “North Africa”, “War Front (WWI)”, “War Front (WWII)” or “Johnson Space Center”.

330 - `people_mentioned`: list of disambiguated people mentioned in the article. Each disambiguated  
331 person is represented as a dictionary, where `wikidata_id` is their ID in Wikidata, `person_name` is  
332 their name on Wikipedia, `person_gender` is their gender from Wikidata and `person_occupation`  
333 is the first listed occupation on Wikidata.

334 - `cluster_size`: Number of newspapers that ran the wire article. Equals length of  
335 `newspaper_metadata`.

336 The whole dataset can be easily downloaded using the `datasets` library:

```
337 from datasets import load_dataset
338 dataset_dict = load_dataset("dell-research-harvard/newswire")
```

339 Specific files can be downloaded by specifying them:

```
340 from datasets import load_dataset
341 load_dataset(
342     "dell-research-harvard/newswire",
343     data_files=["1929_data_clean.json", "1969_data_clean.json"]
344 )
```

### 345 **3.6 Author statement**

346 We bear all responsibility in case of violation of rights.

### 347 **3.7 Hosting, licensing and maintenance Plan**

348 We have chosen to host Newswire on huggingface as this ensures long-term access and preservation  
349 of the dataset.

### 350 **3.8 Dataset documentation and intended uses**

351 We follow the datasheets for datasets template (8).

352 **3.8.1 Motivation**

353 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a  
354 specific gap that needed to be filled? Please provide a description.

355 *The dataset was created to provide researchers with a large, high-quality corpus of structured*  
356 *and transcribed newspaper article texts from American newswires. These texts provide a massive*  
357 *repository of information about historical topics and events. The dataset will be useful to a wide*  
358 *variety of researchers including historians, other social scientists, and NLP practitioners.*

359 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**  
360 **company, institution, organization)?**

361 *NewsWire was created by a team of researchers at Harvard University.*

362 **Who funded the creation of the dataset?** If there is an associated grant, please provide the name  
363 of the grantor and the grant name and number.

364 *The creation of the dataset was funded by the Harvard Data Science Initiative, and the Harvard*  
365 *Economics Department Ken Griffin Fund. Compute credits provided by Microsoft Azure to the*  
366 *Harvard Data Science Initiative.*

367 **Any other comments?**

368 *None.*

369 **3.8.2 Composition**

370 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**  
371 **countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and  
372 interactions between them; nodes and edges)? Please provide a description.

373 *NewsWire comprises instances of newspaper articles. Accompanying each article is a list of newspa-*  
374 *pers that ran the article, classification of whether the article is about certain topics, a list of entities*  
375 *detected in the article, and a disambiguation of people mentioned in the article.*

376 **How many instances are there in total (of each type, if appropriate)?**

377 *NewsWire contains 2,719,607 unique articles.*

378 **Does the dataset contain all possible instances or is it a sample (not necessarily random)**  
379 **of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the  
380 sample representative of the larger set (e.g., geographic coverage)? If so, please describe how  
381 this representativeness was validated/verified. If it is not representative of the larger set, please  
382 describe why not (e.g., to cover a more diverse range of instances, because instances were withheld  
383 or unavailable).

384 *Many newspapers were not preserved, so we cannot guarantee that this dataset contains all possible*  
385 *instances.*

386 **What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or**  
387 **features?** In either case, please provide a description.

388 *Each data instance consists of raw data and derived data. Specifically, an example from NewsWire*  
389 *is:*

```
390     {  
391         "year": 1880,  
392         "dates": ["Feb-23-1880"],  
393         "article": "SENATE Washington, Feb. 23.--Bayard moved that in respect of the
```

```

394         memory of George Washington the senate adjourn ... ",
395     "byline": "",
396     "newspaper_metadata": [
397         {
398             "lccn": "sn92053943",
399             "newspaper_title": "the rock island argus",
400             "newspaper_city": "rock island",
401             "newspaper_state": " illinois "
402         },
403         ...
404     ],
405     "antitrust": 0,
406     "civil_rights": 0,
407     "crime": 0,
408     "govt_regulation": 1,
409     "labor_movement": 0,
410     "politics": 1,
411     "protests": 0,
412     "ca_topic": "Federal Government Operations",
413     "ner_words": ["SENATE", "Washington", "Feb", "23", "Bayard", "moved", "that",
414                 "in", "respect", "of", "the", "memory", "of", "George", "Washington",
415                 "the", "senate", "adjourn", ... ],
416     "ner_labels": ["B-ORG", "B-LOC", "O", "B-PER", "B-PER", "O", "O", "O", "O",
417                  "O", "O", "O", "O", "B-PER", "I-PER", "O", "B-ORG", "O", ...],
418     "wire_city": "Washington",
419     "wire_state": "district of columbia",
420     "wire_country": "United States",
421     "wire_coordinates": [38.89511, -77.03637],
422     "wire_location_notes": "",
423     "people_mentioned": [
424         {
425             "wikidata_id": "Q23",
426             "person_name": "George Washington",
427             "person_gender": "man",
428             "person_occupation": "politician"
429         },
430         ...
431     ],
432     "cluster_size": 8
433 }

```

434 The data fields are:

- 435 - year: year of article publication.
- 436 - dates: list of dates on which this article was published, as strings in the form mmm-DD-YYYY.
- 437 - byline: article byline, if any.
- 438 - article: article text.
- 439 - newspaper\_metadata: list of newspapers that carried the article. Each newspaper is represented as a list of dictionaries, where lccn is the newspaper's Library of Congress identifier, newspaper\_title is the name of the newspaper, and newspaper\_city and newspaper\_state give the location of the newspaper.
- 443 - antitrust: binary variable. 1 if the article was classified as being about antitrust.

444 - `civil_rights`: binary variable. 1 if the article was classified as being about civil rights.

445 - `crime`: binary variable. 1 if the article was classified as being about crime.

446 - `govt_regulation`: binary variable. 1 if the article was classified as being about government  
447 regulation.

448 - `labor_movement`: binary variable. 1 if the article was classified as being about the labor movement.

449 - `politics`: binary variable. 1 if the article was classified as being about politics.

450 - `protests`: binary variable. 1 if the article was classified as being about protests.

451 - `ca_topic`: predicted Comparative Agendas topic of article.

452 - `wire_city`: City of wire service bureau that wrote the article.

453 - `wire_state`: State of wire service bureau that wrote the article.

454 - `wire_country`: Country of wire service bureau that wrote the article.

455 - `wire_coordinates`: Coordinates of city of wire service bureau that wrote the article.

456 - `wire_location_notes`: Contains wire dispatch location if it is not a geographic location. Can  
457 be one of “Pacific Ocean (WWII)”, “Supreme Headquarters Allied Expeditionary Force (WWII)”,  
458 “North Africa”, “War Front (WWI)”, “War Front (WWII)” or “Johnson Space Center”.

459 - `people_mentioned`: list of disambiguated people mentioned in the article. Each disambiguated  
460 person is represented as a dictionary, where `wikidata_id` is their ID in Wikidata, `person_name` is  
461 their name on Wikipedia, `person_gender` is their gender from Wikidata and `person_occupation`  
462 is the first listed occupation on Wikidata.

463 - `cluster_size`: Number of newspapers that ran the wire article. Equals length of  
464 `newspaper_metadata`.

465 **Is there a label or target associated with each instance?** If so, please provide a description.

466 *The data is not labelled, but has had inference from multiple models run on it.*

467 **Is any information missing from individual instances?** If so, please provide a description,  
468 explaining why this information is missing (e.g., because it was unavailable). This does not include  
469 intentionally removed information, but might include, e.g., redacted text.

470 *In some cases, there may be no byline, as the article did not have one, or it was not detected.*  
471 *wire\_city, wire\_state, wire\_country, wire\_coordinates are missing when no location was*  
472 *detected. person\_gender and person\_occupation are missing if no gender or occupation was*  
473 *listed on Wikidata.*

474 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social  
475 network links)?** If so, please describe how these relationships are made explicit.

476 *No relationships between instances are made explicit.*

477 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so,  
478 please provide a description of these splits, explaining the rationale behind them.

479 *There are no recommended splits.*

480 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a  
481 description.

482 *The data is sourced from OCR’d text of historical newspapers. Therefore some of the texts contain*  
483 *OCR errors.*

484 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**  
485 **websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there  
486 guarantees that they will exist, and remain constant, over time; b) are there official archival versions  
487 of the complete dataset (i.e., including the external resources as they existed at the time the dataset  
488 was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external  
489 resources that might apply to a future user? Please provide descriptions of all external resources and  
490 any restrictions associated with them, as well as links or other access points, as appropriate.

491 *The data is self-contained.*

492 **Does the dataset contain data that might be considered confidential (e.g., data that is pro-**  
493 **ected by legal privilege or by doctor-patient confidentiality, data that includes the content of**  
494 **individuals non-public communications)?** If so, please provide a description.

495 *The dataset does not contain information that might be viewed as confidential.*

496 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,**  
497 **or might otherwise cause anxiety?** If so, please describe why.

498 *The headlines in the dataset reflect diverse attitudes and values from the period in which they were*  
499 *written, 1878-1977, and contain content that may be considered offensive for a variety of reasons.*

500 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

501 *Many news articles are about people.*

502 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how  
503 these subpopulations are identified and provide a description of their respective distributions within  
504 the dataset.

505 *The dataset does not specifically identify any subpopulations.*

506 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or**  
507 **indirectly (i.e., in combination with other data) from the dataset?** If so, please describe how.

508 *If an individual appeared in the news during this period, then article text may contain their name,*  
509 *age, and information about their actions. Further, for prominent individuals, we have disambiguated*  
510 *them to Wikipedia, which directly identifies them.*

511 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that**  
512 **reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or**  
513 **union memberships, or locations; financial or health data; biometric or genetic data; forms of**  
514 **government identification, such as social security numbers; criminal history)?** If so, please  
515 provide a description.

516 *All information that it contains is already publicly available in the newspapers used to create the*  
517 *data.*

518 **Any other comments?**

519 *None.*

### 520 **3.8.3 Collection Process**

521 **How was the data associated with each instance acquired?** Was the data directly observable (e.g.,  
522 raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived  
523 from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If data was  
524 reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If  
525 so, please describe how.

526 *The dataset combines raw data and derived data. The pipeline used to extract the data and to create*  
527 *the derived data is described in detail within the paper. The dataset described here is the output of*  
528 *that pipeline.*

529 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sen-**  
530 **sor, manual human curation, software program, software API)?** How were these mechanisms  
531 or procedures validated?

532 *These methods are described in detail in the main text and supplementary materials of this paper.*

533 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**  
534 **probabilistic with specific sampling probabilities)?**

535 *The dataset was not sampled from a larger set.*

536 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and**  
537 **how were they compensated (e.g., how much were crowdworkers paid)?**

538 *We used student annotators to create the validation and test sets. They were paid \$15 per hour, a rate*  
539 *set by a Harvard economics department program providing research assistantships for undergradu-*  
540 *ates.*

541 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe**  
542 **of the data associated with the instances (e.g., recent crawl of old news articles)?** If not, please  
543 describe the timeframe in which the data associated with the instances was created.

544 *The articles were written between 1878 and 1977. They were processed between 2020 and 2024.*

545 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so,  
546 please provide a description of these review processes, including the outcomes, as well as a link or  
547 other access point to any supporting documentation.

548 *No, this dataset uses entirely public information and hence does not fall under the domain of*  
549 *Harvard's institutional review board.*

550 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

551 *Historical newspapers contain a variety of information about people.*

552 **Did you collect the data from the individuals in question directly, or obtain it via third parties**  
553 **or other sources (e.g., websites)?**

554 *The data were obtained from historical newspapers.*

555 **Were the individuals in question notified about the data collection?** If so, please describe (or  
556 show with screenshots or other information) how notice was provided, and provide a link or other  
557 access point to, or otherwise reproduce, the exact language of the notification itself.

558 *Individuals were not notified; the data came from publicly available newspapers.*

559 **Did the individuals in question consent to the collection and use of their data?** If so, please  
560 describe (or show with screenshots or other information) how consent was requested and provided,  
561 and provide a link or other access point to, or otherwise reproduce, the exact language to which the  
562 individuals consented.

563 *The dataset was created from publicly available historical newspapers.*

564 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke**  
565 **their consent in the future or for certain uses?** If so, please provide a description, as well as a  
566 link or other access point to the mechanism (if appropriate).

567 *Not applicable.*

568 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data**  
569 **protection impact analysis) been conducted?** If so, please provide a description of this analysis,  
570 including the outcomes, as well as a link or other access point to any supporting documentation.

571 *No.*

572 **Any other comments?**

573 *None.*

#### 574 **3.8.4 Preprocessing/cleaning/labeling**

575 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**  
576 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**  
577 **of missing values)?** If so, please provide a description. If not, you may skip the remainder of the  
578 questions in this section.

579 *See the description in the main text.*

580 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support**  
581 **unanticipated future uses)?** If so, please provide a link or other access point to the “raw” data.

582 *All data is in the dataset.*

583 **Is the software used to preprocess/clean/label the instances available?** If so, please provide a  
584 link or other access point.

585 *No specific software was used to clean the instances.*

586 **Any other comments?**

587 *None.*

#### 588 **3.8.5 Uses**

589 **Has the dataset been used for any tasks already?** If so, please provide a description.

590 *No.*

591 **Is there a repository that links to any or all papers or systems that use the dataset?** If so,  
592 please provide a link or other access point.

593 *No such repository currently exists.*

594 **What (other) tasks could the dataset be used for?**

595 *There are a large number of potential uses in the social sciences, digital humanities, and deep*  
596 *learning research, discussed in more detail in the main text.*

597 **Is there anything about the composition of the dataset or the way it was collected and prepro-**  
598 **cessed/cleaned/labeled that might impact future uses?** For example, is there anything that a  
599 future user might need to know to avoid uses that could result in unfair treatment of individuals or  
600 groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms,  
601 legal risks) If so, please provide a description. Is there anything a future user could do to mitigate  
602 these undesirable harms?

603 *This dataset contains unfiltered content composed by newspaper editors, columnists, and other*  
604 *sources. It reflects their biases and any factual errors that they made.*

605 **Are there tasks for which the dataset should not be used?** If so, please provide a description.

606 *We would urge caution in using the data to train generative language models - without additional*  
607 *filtering - as it contains content that many would consider toxic.*

608 **Any other comments?**

609 *None*

### 610 **3.8.6 Distribution**

611 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**  
612 **organization) on behalf of which the dataset was created?** If so, please provide a description.

613 *Yes. The dataset is available for public use.*

614 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)** Does the dataset  
615 have a digital object identifier (DOI)?

616 *The dataset is hosted on huggingface. Its DOI is 10.57967/hf/2423.*

617 **When will the dataset be distributed?**

618 *The dataset was distributed on 7th June 2024.*

619 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**  
620 **and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and  
621 provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU,  
622 as well as any fees associated with these restrictions.

623 *The dataset is distributed under a Creative Commons CC-BY license. The terms of this license can be*  
624 *viewed at <https://creativecommons.org/licenses/by/2.0/>*

625 **Have any third parties imposed IP-based or other restrictions on the data associated with**  
626 **the instances?** If so, please describe these restrictions, and provide a link or other access point  
627 to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these  
628 restrictions.

629 *There are no third party IP-based or other restrictions on the data.*

630 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**  
631 **instances?** If so, please describe these restrictions, and provide a link or other access point to, or  
632 otherwise reproduce, any supporting documentation.

633 *No export controls or other regulatory restrictions apply to the dataset or to individual instances.*

634 **Any other comments?**

635 *None.*

### 636 **3.8.7 Maintenance**

637 **Who will be supporting/hosting/maintaining the dataset?**

638

639 *The dataset is hosted on huggingface.*

640 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

641



642 *The recommended method of contact is using the huggingface ‘community’ capacity. Additionally,*  
643 *Melissa Dell can be contacted at melissadell@fas.harvard.edu.*

644 **Is there an erratum?** If so, please provide a link or other access point.

645 *There is no erratum.*

646 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**

647 If so, please describe how often, by whom, and how updates will be communicated to users (e.g.,  
648 mailing list, GitHub)?

649 *If we update the dataset, we will notify users via the huggingface Dataset Card.*

650 **If the dataset relates to people, are there applicable limits on the retention of the data associated**  
651 **with the instances (e.g., were individuals in question told that their data would be retained for a**  
652 **fixed period of time and then deleted)?** If so, please describe these limits and explain how they  
653 will be enforced.

654 *There are no applicable limits on the retention of data.*

655 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please  
656 describe how. If not, please describe how its obsolescence will be communicated to users.

657 *If we update the dataset, older versions of the dataset will not continue to be hosted. We will notify*  
658 *users via the huggingface Dataset Card.*

659 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**  
660 **them to do so?** If so, please provide a description. Will these contributions be validated/verified?  
661 If so, please describe how. If not, why not? Is there a process for communicating/distributing these  
662 contributions to other users? If so, please provide a description.

663 *Others can contribute to the dataset using the huggingface ‘community’ capacity. This allows for*  
664 *anyone to ask questions, make comments and submit pull requests. We will validate these pull requests.*  
665 *A record of public contributions will be maintained on huggingface, allowing communication to other*  
666 *users.*

667 **Any other comments?**

668 *None.*

## 669 **References**

670 [1] BIEWALD, L. Experiment tracking with weights and biases, 2020. Software available from  
671 wandb.com.

672 [2] BOYDSTUN, A. E. *Making the news : politics, the media, and agenda setting.* The University  
673 of Chicago Press, Chicago ; London, 2013.

674 [3] BRYAN, T., CARLSON, J., ARORA, A., AND DELL, M. Efficientocr: An extensible, open-  
675 source package for efficiently digitizing world knowledge”. *Empirical Methods on Natural*  
676 *Language Processing (Systems Demonstrations Track) (2023).*

677 [4] BUITINCK, L., LOUPPE, G., BLONDEL, M., PEDREGOSA, F., MUELLER, A., GRISEL, O.,  
678 NICULAE, V., PRETTENHOFER, P., GRAMFORT, A., GROBLER, J., LAYTON, R., VANDER-  
679 PLAS, J., JOLY, A., HOLT, B., AND VAROQUAUX, G. API design for machine learning  
680 software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for*  
681 *Data Mining and Machine Learning (2013)*, pp. 108–122.

- 682 [5] CARLSON, J., BRYAN, T., AND DELL, M. Efficient ocr for building a diverse digital history.  
683 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*  
684 (forthcoming).
- 685 [6] DELL, M., CARLSON, J., BRYAN, T., SILCOCK, E., ARORA, A., SHEN, Z., D’AMICO-  
686 WONG, L., LE, Q., QUERUBIN, P., AND HELDRING, L. American stories: A large-scale  
687 structured text dataset of historical us newspapers. *Advances in Neural Information and*  
688 *Processing Systems, Datasets and Benchmarks* (2023).
- 689 [7] FRANKLIN, B., SILCOCK, E., ARORA, A., BRYAN, T., AND DELL, M. News deja vu:  
690 Connecting past and present with semantic search, 2024.
- 691 [8] GEBRU, T., MORGENSTERN, J., VECCHIONE, B., VAUGHAN, J. W., WALLACH, H., AU2, H.  
692 D. I., AND CRAWFORD, K. Datasheets for datasets, 2021.
- 693 [9] HADSELL, R., CHOPRA, S., AND LECUN, Y. Dimensionality reduction by learning an  
694 invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and*  
695 *Pattern Recognition (CVPR’06)* (2006), vol. 2, IEEE, pp. 1735–1742.
- 696 [10] HOWARD, A., SANDLER, M., CHU, G., CHEN, L.-C., CHEN, B., TAN, M., WANG, W.,  
697 ZHU, Y., PANG, R., VASUDEVAN, V., ET AL. Searching for mobilenetv3. In *Proceedings of*  
698 *the IEEE/CVF international conference on computer vision* (2019), pp. 1314–1324.
- 699 [11] JOHNSON, J., DOUZE, M., AND JÉGOU, H. Billion-scale similarity search with gpus. *IEEE*  
700 *Transactions on Big Data* 7, 3 (2019), 535–547.
- 701 [12] KARPUKHIN, V., OĞUZ, B., MIN, S., LEWIS, P., WU, L., EDUNOV, S., CHEN, D., AND  
702 YIH, W.-T. Dense passage retrieval for open-domain question answering. *arXiv preprint*  
703 *arXiv:2004.04906* (2020).
- 704 [13] LEWIS, M., LIU, Y., GOYAL, N., GHAZVININEJAD, M., MOHAMED, A., LEVY, O., STOY-  
705 ANOV, V., AND ZETTLEMOYER, L. Bart: Denoising sequence-to-sequence pre-training for  
706 natural language generation, translation, and comprehension, 2019.
- 707 [14] LIU, Y., OTT, M., GOYAL, N., DU, J., JOSHI, M., CHEN, D., LEVY, O., LEWIS, M.,  
708 ZETTLEMOYER, L., AND STOYANOV, V. Roberta: A robustly optimized bert pretraining  
709 approach. *arXiv preprint arXiv:1907.11692* (2019).
- 710 [15] PASZKE, A., GROSS, S., MASSA, F., LERER, A., BRADBURY, J., CHANAN, G., KILLEEN,  
711 T., LIN, Z., GIMELSHEIN, N., ANTIGA, L., DESMAISON, A., KOPF, A., YANG, E., DEVITO,  
712 Z., RAISON, M., TEJANI, A., CHILAMKURTHY, S., STEINER, B., FANG, L., BAI, J., AND  
713 CHINTALA, S. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In  
714 *Advances in Neural Information Processing Systems 32* (2019), H. Wallach, H. Larochelle,  
715 A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, Eds., Curran Associates, Inc., pp. 8024–  
716 8035.
- 717 [16] REIMERS, N., AND GUREVYCH, I. Sentence-bert: Sentence embeddings using siamese  
718 bert-networks. *arXiv preprint arXiv:1908.10084* (2019).
- 719 [17] SILCOCK, E., ARORA, A., AND DELL, M. A massive scale semantic similarity dataset of  
720 historical english. *Advances in Neural Information and Processing Systems, Datasets and*  
721 *Benchmarks* (2023).
- 722 [18] SILCOCK, E., D’AMICO-WONG, L., YANG, J., AND DELL, M. Noise-robust de-duplication  
723 at scale. In *The Eleventh International Conference on Learning Representations* (2022).
- 724 [19] SONG, K., TAN, X., QIN, T., LU, J., AND LIU, T.-Y. Mpnet: Masked and permuted pre-  
725 training for language understanding. *Advances in Neural Information Processing Systems 33*  
726 (2020), 16857–16867.

- 727 [20] ULTALYTICS. Yolo v8 github repository. [https://github.com/ultralytics/](https://github.com/ultralytics/ultralytics)  
728 [ultralytics](https://github.com/ultralytics/ultralytics), 2023.
- 729 [21] VON PLATEN, P. Transformer-based encoder-decoder models, 2020.
- 730 [22] WIGHTMAN, R. Pytorch image models. [https://github.com/rwightman/](https://github.com/rwightman/pytorch-image-models)  
731 [pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 732 [23] WILLIAMS, A., NANGIA, N., AND BOWMAN, S. A broad-coverage challenge corpus for  
733 sentence understanding through inference. In *Proceedings of the 2018 Conference of the*  
734 *North American Chapter of the Association for Computational Linguistics: Human Language*  
735 *Technologies, Volume 1 (Long Papers)* (2018), Association for Computational Linguistics,  
736 pp. 1112–1122.
- 737 [24] WU, L., PETRONI, F., JOSIFOSKI, M., RIEDEL, S., AND ZETTLEMOYER, L. Scalable  
738 zero-shot entity linking with dense entity retrieval. *arXiv preprint arXiv:1911.03814* (2019).