

1 **Supplementary Materials**

2 **S-1 Methods to Associate Articles**

3 Figure 6 in the main text illustrates the full article association procedure.

4 First, we used a rule-based algorithm using associate article bounding boxes that are under the same
5 headline, as these are part of the same article with extremely high probability. Algorithm 1 gives
6 pseudocode for this method. We set the parameters as $P_S = 100$, $P_T = 20$, $P_B = 50$.

7 For training data, where we want article pairs that are not only part of the same article, but also where
8 they appear in the given order, we further narrow down the pairs. Specifically, we use only those pairs
9 which are horizontally next to each other, and which have no other bounding boxes below them, as
10 for these pairs, we can guarantee that the pair of bounding follow directly after one another (whereas
11 for other article bounding boxes that share a headline, there may be a third bounding box in between).
12 Algorithm 2 shows pseudocode for this procedure, and we used $P_C = 5$, and it is further illustrated
13 in panel A of figure 6 in the main text.

14 For hard negatives, we used article boxes under the same headline in reverse reading order (right to
15 left). For standard negatives, we took pairs of articles on the same page, where B was above and to
16 the left of A, as articles do not read from right to left. One twelfth of our training data were positive
17 pairs, another twelfth were hard negative pairs and the remainder were standard negative pairs. This
18 outperformed a more balanced training sample.

19 We use this dataset to finetune a cross-encoder using a RoBERTa base model [4]. We used a Bayesian
20 search algorithm [1] to find optimal hyperparameters on one tenth of our training data (limited
21 compute prevented us from running this search with the full dataset), which led to a learning rate
22 of $1.7e-5$, with a batch size of 64 and 29.57% warm up. We trained for 26 epochs with an AdamW
23 optimizer, and optimize a binary cross-entropy loss.

24 We evaluate these methods on a hand-labeled dataset of 214 scans, randomly selected from 1968
25 and 1955. These scans were labeled by a highly-trained undergraduate research assistant. Summary
statistics of this dataset are given in table S-1 and evaluation results are given in the main text.

Scan count	Article bounding boxes	Headline bounding boxes	Article-article associations
214	3,803	2,805	1,851

Table S-1: Descriptive statistics of article association training data.

26

27 **S-2 Methods to Detect Reproduced Content**

28 To detect reproduced content, we use the contrastively trained bi-encoder model developed by [6],
29 which is trained to learn similar representations for reproduced articles and dissimilar representations
30 for non-reproduced articles. This model is based on an S-BERT MPNET model [5, 7] and is fine-
31 tuned on a hand-labelled dataset of articles from the same underlying wire source, using S-BERT's
32 online contrastive loss [3] implementation, with a 0.2 margin and cosine similarity as the distance
33 metric. The learning rate is $2e-5$ with 100% warm up and a batch size of 32. It uses an AdamW
34 optimizer, and the model is trained for 16 epochs. This bi-encoder is trained and evaluated on a
35 hand-labeled dataset, which is detailed in S-2. The results of this evaluation are given in the main
36 text.

37 To create clusters from the bi-encoder embeddings, we use highly scalable single-linkage clustering,
38 with a cosine similarity threshold of 0.94. We build a graph using articles as nodes, and add edges if
39 the cosine similarity is above this threshold. As edge weights we use the negative exponential of the
40 difference in dates (in days) between the two articles. We then apply Leiden community detection to
41 the graph to control false positive edges that can otherwise merge disparate groups of articles.

Algorithm 1 Rule-based association of article bounding boxes

INPUT: $b_1, \dots, b_n \in B$: set of bounding boxes that appear on the same scan, with their coordinates, denoted $left(b_i)$, $right(b_i)$, $top(b_i)$, $bottom(b_i)$, and type (headline, article, byline etc.), denoted $type(b_i)$.

PARAMETERS:

W : width of scan

H : height of scan

P_S : fraction of width, for creating side margin

P_T : fraction of height, for creating top margin

P_B : fraction of height, for creating bottom margin

OUTPUT: $ArticleArticlePairs = \{(b_i, b_j) \in B \times B \mid b_i, b_j \text{ predicted to be part of the same full article and } type(b_i) = type(b_j) = \text{article}\}$

- 1: Initialise: $M_S = W/P_S$, the side margin, $M_B = H/P_B$, the bottom margin, $M_T = H/P_T$, the top margin, $MatchedHeadlines = \{\}$, $HeadlineArticlePairs = \{\}$, $ArticleArticlePairs = \{\}$
 - 2: **for** all b_0 in B where $type(b_0)$ is article **do**
 - 3: Create $B_0 \subset B$ where:
 - a. All bounding boxes of type byline are removed
 - b. b_0 is removed
 - c. All bounding boxes are removed that do not share at least M_S of the horizontal axis
 - d. All bounding boxes are removed whose bottom is more than M_B below the top of b_0
 - e. All bounding boxes are removed whose bottom is more than M_T above the top of b_0
 - 4: **if** B_0 is not empty **then**
 - 5: Let b_1 be the element of B_0 that has the lowest bottom coordinate
 - 6: **if** $type(b_1)$ is headline **then**
 - 7: $MatchedHeadlines = MatchedHeadlines \cup \{b_1\}$
 - 8: $HeadlineArticlePairs = HeadlineArticlePairs \cup \{(b_0, b_1)\}$
 - 9: **end if**
 - 10: **end if**
 - 11: **end for**
 - 12: **for** all b_h in $MatchedHeadlines$ **do**
 - 13: Let $H_1 \subset HeadlineArticlePairs$ be all pairs that contain that headline, b_h
 - 14: **if** H_1 has at least two elements **then**
 - 15: Let A be all the bounding boxes of type article from the pairs in H_1
 - 16: Let C be all combinations of 2 elements of A
 - 17: $ArticleArticlePairs = ArticleArticlePairs \cup C$
 - 18: **end if**
 - 19: **end for**
-

Algorithm 2 Selection of ordered article pairs

INPUT: $ArticleArticlePairs$, B , from algorithm 1.

PARAMETERS:

P_C : fraction of column width, for creating margin

OUTPUT: $OrderedPairs \subset ArticleArticlePairs$

- 1: Initialise: $OrderedPairs = \{\}$
 - 2: **for** p in $ArticleArticlePairs$ **do**
 - 3: Let p_l be the element of p with the furthest left coordinate
 - 4: Let p_r be the other element
 - 5: **if** $left(p_r)$ is not further to the right of $right(p_l)$ than $width(p_l)/P_C$ **then**
 - 6: **if** there are no other bounding boxes below p_l **then**
 - 7: $OrderedPairs = OrderedPairs \cup \{p\}$
 - 8: **end if**
 - 9: **end if**
 - 10: **end for**
-

	Positives Pairs	Negative Pairs	Reproduced Articles	Singleton Articles	Total Articles
Training Data					
Training	36,291	37,637	891	–	7,728
Validation	3,042	3,246	20	–	283
Full Day Evaluation					
Validation	28,547	12,409,031	447	2,162	4,988
Test	54,996	100,914,159	1,236	8,046	14,211
Full Dataset	122,876	113,364,073	2,594	10,208	27,210

Table S-2: Summary statistics of training and evaluation data for detecting duplicate content.

42 We further remove clusters that have over 50 articles and contain articles with greater than five
43 different dates. We also remove clusters that contain over 50 articles, when the number of articles
44 is more than double the number of unique newspapers from which these articles are sourced. This
45 removes clusters of content that are correctly clustered in the sense of being based on the same
46 underlying source, but are not useful for the HEADLINES dataset. For example, an advertisement
47 (misclassified as an article due to an article-like appearance) might be repeated by the same newspaper
48 on multiple different dates and would be removed by these rules, or weather forecasts can be very
49 near duplicates across space and time, forming large clusters.

50 S-3 A Summary of Copyright Law for Works Published in the United States

Date of Publication	Conditions	Copyright Term
Public Domain		
Anytime	Works prepared by an officer/employee of the U.S. Government as part of their official duties	None
Before 1928	None	None. Copyright expired.
1928 through 1977	Published without a copyright notice	None. Failure to comply with required formalities
1978 to 1 March 1989	Published without notice and without subsequent registration within 5 years	None. Failure to comply with required formalities
1928 through 1963	Published with notice but copyright was not renewed	None. Copyright expired
Copyrighted		
1978 to 1 March 1989	Published without notice, but with subsequent registration within 5 years	70 (95) years after the death of author (corporate author)
1928 through 1963	Published with notice and the copyright was renewed	95 years after publication
1964 through 1977	Published with notice	95 years after publication
1978 to 1 March 1989	Created after 1977 and published with notice	70 (95) years after the death of author (corporate author) or 120 years after creation, if earlier
1978 to 1 March 1989	Created before 1978 and first published with notice in the specified period	The greater of the term specified in the previous entry or 31 December 2047
From 1 March 1989 through 2002	Created after 1977	70 (95) years after the death of author (corporate author) or 120 years after creation, if earlier
From 1 March 1989 through 2002	Created before 1978 and first published in this period	The greater of the term specified in the previous entry or 31 December 2047
After 2002	None	70 (95) years after the death of author (corporate author) or 120 years after creation, if earlier

Table S-3: This table summarizes U.S. copyright law, based on a similar table produced by the Cornell libraries. For concision, we focus on works initially published in the United States. A variety of other cases are also covered at <https://guides.library.cornell.edu/copyright>.

51 **S-4 Dataset Information**

52 **S-4.1 Dataset URL**

53 HEADLINES can be found at [https://huggingface.co/datasets/dell-research-harvard](https://huggingface.co/datasets/dell-research-harvard/headlines-semantic-similarity)
54 [/headlines-semantic-similarity](https://huggingface.co/datasets/dell-research-harvard/headlines-semantic-similarity).

55 This dataset has structured metadata following schema.org, and is readily discoverable.¹

56 **S-4.2 DOI**

57 The DOI for this dataset is: [10.57967/hf/0751](https://doi.org/10.57967/hf/0751).

58 **S-4.3 License**

59 HEADLINES has a Creative Commons CC-BY license.

60 **S-4.4 Dataset usage**

61 The dataset is hosted on huggingface, in json format. Each year in the dataset is divided into a distinct
62 file (eg. `1952_headlines.json`).

63 The data is presented in the form of clusters, rather than pairs to eliminate duplication of text data
64 and minimize the storage size of the datasets.

65 An example from HEADLINES looks like:

```
66 {  
67     "headline": "FRENCH AND BRITISH BATTLESHIPS IN MEXICAN WATERS",  
68     "group_id": 4  
69     "date": "May-14-1920",  
70     "state": "kansas",  
71 }
```

72 The data fields are:

- 73 • `headline`: headline text.
- 74 • `date`: the date of publication of the newspaper article, as a string in the form `mmm-DD-`
75 `YYYY`.
- 76 • `state`: state of the newspaper that published the headline.
- 77 • `group_id`: a number that is shared with all other headlines for the same article. This number
78 is unique across all year files.

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

```
80 from datasets import load_dataset  
81 dataset_dict = load_dataset("dell-research-harvard/headlines-semantic-similarity")
```

82 Specific files can be downloaded by specifying them:

```
83 from datasets import load_dataset  
84 load_dataset(  
85     "dell-research-harvard/headlines-semantic-similarity",  
86     data_files=["1929_headlines.json", "1989_headlines.json"]  
87 )
```

¹See https://search.google.com/test/rich-results/result?id=_HKjxIv-LaF_8E1AarsM_g for full metadata.

88 **S-4.5 Author statement**

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

90 **S-4.6 Maintenance Plan**

91 We have chosen to host HEADLINES on huggingface as this ensures long-term access and preservation
92 of the dataset.

93 **S-4.7 Dataset documentation and intended uses**

94 We follow the datasheets for datasets template [2].

95 **S-4.7.1 Motivation**

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

98 *Transformer language models contrastively trained on large-scale semantic similarity datasets are*
99 *integral to a variety of applications in natural language processing (NLP). A variety of semantic*
100 *similarity datasets have been used for this purpose, with positive text pairs related to each other*
101 *in some way. Many of these datasets are relatively small, and the bulk of the larger datasets are*
102 *created from recent web texts; e.g. positives are drawn from the texts in an online comment thread or*
103 *duplicate questions in a forum. Relative to existing datasets, HEADLINES is very large, covering a*
104 *vast array of topics. This makes it useful generally speaking for semantic similarity pre-training. It*
105 *also covers a long period of time, making it a rich training data source for the study of historical*
106 *texts and semantic change. It captures semantic similarity directly, as the positive pairs summarize*
107 *the same underlying texts.*

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

110 *HEADLINES was created by Melissa Dell and Emily Silcock, at Harvard University.*

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

113 *The creation of the dataset was funded by the Harvard Data Science Initiative, Harvard Catalyst, and*
114 *compute credits provided by Microsoft Azure to the Harvard Data Science Initiative.*

115 **Any other comments?**

116 *None.*

117 **S-4.7.2 Composition**

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

121 *HEADLINES comprises instances of newspaper headlines and relationships between them. Specifically,*
122 *each headline includes information on the text of the headline, the date of publication, and the state*
123 *it was published in. Headlines have relationships between them if they are semantic similarity pairs,*
124 *that is, if they two different headlines for the same newspaper article.*

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

126 *HEADLINES contains 34,867,488 different headlines and 396,001,930 positive relationships between*
127 *headlines.*

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

134 *Many local newspapers were not preserved, and newspapers with the widest circulation tended to*
135 *renew their copyrights, so cannot be included.*

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

138 *Each data instance consists of raw data. Specifically, an example from HEADLINES is:*

```
139 {  
140   "headline": "FRENCH AND BRITISH BATTLESHIPS IN MEXICAN WATERS",  
141   "group_id": 4  
142   "date": "May-14-1920",  
143   "state": "kansas",  
144 }
```

145 *The data fields are:*

- 146 • *headline: headline text.*
- 147 • *date: the date of publication of the newspaper article, as a string in the form mmm-DD-*
148 *YYYY.*
- 149 • *state: state of the newspaper that published the headline.*
- 150 • *group_id: a number that is shared with all other headlines for the same article. This*
151 *number is unique across all year files.*

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

153 *Each instance contains a group_id as mentioned directly above. This is a number that is shared by*
154 *all other instances that are positive semantic similarity pairs.*

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

158 *In some cases, the state of publication is missing, due to incomplete metadata.*

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

161 *Relationships between instances are made explicit in the group_id variable, as detailed above.*

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

164 *There are no recommended splits.*

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

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

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

176 *The data is self-contained.*

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

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

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

183 *The headlines in the dataset reflect diverse attitudes and values from the period in which they were*
184 *written, 1920-1989, and contain content that may be considered offensive for a variety of reasons.*

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

186 *Many news articles are about people.*

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

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

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

193 *If an individual appeared in the news during this period, then headline text may contain their name,*
194 *age, and information about their actions.*

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

200 *All information that it contains is already publicly available in the newspapers used to create the*
201 *headline pairs.*

202 **Any other comments?**

203 *None.*

204 **S-4.7.3 Collection Process**

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

210 *To create HEADLINES, we digitized front pages from off-copyright newspapers spanning 1920-1989.*
211 *Historically, around half of articles in U.S. local newspapers came from newswires like the Associated*
212 *Press. While local papers reproduced articles from the newswire, they wrote their own headlines,*
213 *which form abstractive summaries of the associated articles. We associate articles and their headlines*
214 *by exploiting document layouts and language understanding. We then use deep neural methods to*
215 *detect which articles are from the same underlying source, in the presence of substantial noise and*
216 *abridgement. The headlines of reproduced articles form positive semantic similarity pairs.*

217 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)?** How were these mechanisms
218 or procedures validated?
219

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

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

222 *The dataset was not sampled from a larger set.*
223

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

225 *We used student annotators to create the validation sets for associating bounding boxes, and the*
226 *training and validation sets for clustering duplicated articles. They were paid \$15 per hour, a rate set*
227 *by a Harvard economics department program providing research assistantships for undergraduates.*
228

229 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?** If not, please

230 describe the timeframe in which the data associated with the instances was created.
231

232 *The headlines were written between 1920 and 1989. Semantic similarity pairs were computed in*
233 *2023.*

234 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so, please provide a description of these review processes, including the outcomes, as well as a link or
235 other access point to any supporting documentation.
236

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

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

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

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

242 *The data were obtained from off-copyright historical newspapers.*
243

244 **Were the individuals in question notified about the data collection?** If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other
245 access point to, or otherwise reproduce, the exact language of the notification itself.
246

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

248 **Did the individuals in question consent to the collection and use of their data?** If so, please describe (or show with screenshots or other information) how consent was requested and provided,
249 and provide a link or other access point to, or otherwise reproduce, the exact language to which the
250 individuals consented.
251

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

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

256 *Not applicable.*

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

260 *No.*

261 **Any other comments?**

262 *None.*

263 **S-4.7.4 Preprocessing/cleaning/labeling**

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

268 *See the description in the main text.*

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

271 *No.*

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

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

275 **Any other comments?**

276 *None.*

277 **S-4.7.5 Uses**

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

279 *No.*

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

282 *No.*

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

284 *The dataset can be used for training models for semantic similarity, studying language change over*
285 *time and studying difference in language across space.*

286 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
287 **cessed/cleaned/labeled that might impact future uses?** For example, is there anything that a
288 future user might need to know to avoid uses that could result in unfair treatment of individuals or
289 groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms,

290 legal risks) If so, please provide a description. Is there anything a future user could do to mitigate
291 these undesirable harms?

292 *The dataset contains historical news headlines, which will reflect current affairs and events of the*
293 *time period in which they were created, 1920-1989, as well as the biases of this period.*

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

295 *It is intended for training semantic similarity models and studying semantic variation across space*
296 *and time.*

297 **Any other comments?**

298 *None*

299 **S-4.7.6 Distribution**

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

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

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

305 *The dataset is hosted on huggingface. Its DOI is 10.57967/hf/0751.*

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

307 *The dataset was distributed on 7th June 2023.*

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

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

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

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

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

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

323 **Any other comments?**

324 *None.*

325 **S-4.7.7 Maintenance**

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

327

328 *The dataset is hosted on huggingface.*

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

330

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

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

334 *There is no erratum.*

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

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

338 *We have no plans to update the dataset. If we do, we will notify users via the huggingface Dataset*
339 *Card.*

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

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

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

347 *We have no plans to update the dataset. If we do, older versions of the dataset will not continue to be*
348 *hosted. We will notify users via the huggingface Dataset Card.*

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

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

357 **Any other comments?**

358 *None.*

359 **References**

- 360 [1] FALKNER, S., KLEIN, A., AND HUTTER, F. BOHB: Robust and efficient hyperparameter
361 optimization at scale. In *Proceedings of the 35th International Conference on Machine Learning*
362 (10–15 Jul 2018), J. Dy and A. Krause, Eds., vol. 80 of *Proceedings of Machine Learning*
363 *Research*, PMLR, pp. 1437–1446.
- 364 [2] GEBRU, T., MORGENSTERN, J., VECCHIONE, B., VAUGHAN, J. W., WALLACH, H., AU2, H.
365 D. I., AND CRAWFORD, K. Datasheets for datasets, 2021.
- 366 [3] HADSELL, R., CHOPRA, S., AND LECUN, Y. Dimensionality reduction by learning an in-
367 variant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern*
368 *Recognition (CVPR'06)* (2006), vol. 2, IEEE, pp. 1735–1742.
- 369 [4] LIU, Y., OTT, M., GOYAL, N., DU, J., JOSHI, M., CHEN, D., LEVY, O., LEWIS, M.,
370 ZETTLEMOYER, L., AND STOYANOV, V. Roberta: A robustly optimized bert pretraining
371 approach. *arXiv preprint arXiv:1907.11692* (2019).
- 372 [5] REIMERS, N., AND GUREVYCH, I. Sentence-bert: Sentence embeddings using siamese bert-
373 networks. *arXiv preprint arXiv:1908.10084* (2019).
- 374 [6] SILCOCK, E., D'AMICO-WONG, L., YANG, J., AND DELL, M. Noise-robust de-duplication at
375 scale. *International Conference on Learning Representations* (2023).
- 376 [7] SONG, K., TAN, X., QIN, T., LU, J., AND LIU, T.-Y. Mpnet: Masked and permuted pre-training
377 for language understanding. *Advances in Neural Information Processing Systems 33* (2020),
378 16857–16867.