

POLITICS: Pretraining with Same-story Article Comparison for Ideology Prediction and Stance Detection

Anonymous ACL submission

Abstract

Ideology is at the core of political science. Yet, there still does not exist general-purpose tools that can characterize and predict ideology across different genres of text. To this end, we study the training of PLMs using novel ideology-driven pretraining objectives that rely on the comparison of articles that are on the same stories but written by media of different ideologies. We further collect a large-scale dataset consisting of more than 3.6M political news articles for experiments. Our model POLITICS and its variants outperform strong baselines on 10 out of the 11 ideology prediction and stance detection tasks. Our analysis further shows that POLITICS is especially good at understanding long or formally written texts, and is also robust in few-shot learning scenarios.

1 Introduction

Ideology is an ubiquitous factor in political science, journalism, and media studies (Mullins, 1972; Freeden, 2006; Martin, 2015). Decades of work has gone into measuring ideology based on voting data (Poole and Rosenthal, 1985; Lewis et al., 2021), survey results (Preoțiuc-Pietro et al., 2017; Ansolabehere et al., 2008; Kim and Fording, 1998; Gabel and Huber, 2000), social networks (Barberá et al., 2015), campaign donation records (Bonica, 2013), and textual data (Laver et al., 2003; Diermeier et al., 2012a; Gentzkow et al., 2019; Volkens et al., 2021). Each of these approaches has its strengths and weaknesses. For instance, many political figures do not have a voting record, while surveys are expensive and politicians are often unwilling to disclose ideology. By contrast, political text is abundant and ubiquitous. However, language is complex in nature, often domain-specific, and generally unlabeled, making it challenging to work with. There thus remains a strong need for general-purpose tools for measuring ideology using text that can be applied across multiple genres.

News Story: Donald Trump tests positive for COVID-19.

Daily Kos (left): It’s now clear that Donald Trump **lied** to the nation about when he received a positive test for COVID-19. . . . they’re continuing to act as if nothing has changed—and that **disregarding science** and **lying** to the public are the only possible strategies.

The Washington Times (right): *Trump says he’s “doing very well” . . . President Trump thanked the nation for supporting him* Friday night as he left the White House to be hospitalized for COVID-19. *“I want to thank everybody for the tremendous support. . . .”* Mr. Trump said in a video recorded at the White House.

Breitbart (right): *President Donald Trump thanked Americans for their support* on Friday as he traveled to Walter Reed Military Hospital for further care after he was diagnosed with coronavirus. *“I think I’m doing very well. . . .”* Trump said in a video filmed at the White House and posted to social media.

Figure 1: Article snippets by different media on the same news story. Contents that are used to show stances and ideological leanings are highlighted in **bold** (for phrases) and in *italics* (for facts).

Using text as data, computational models for ideology measurement have rapidly expanded and diversified, including statistical methods such as ideal point estimation (Grosz et al., 1999; Shor and McCarty, 2011) and regression (Peterson and Spirling, 2018); probabilistic models such as Naive Bayes (Evans et al., 2007), support vector machines (Yu et al., 2008), and latent variable models (Barberá et al., 2015); and more recent neural architectures such as recurrent neural networks (Iyyer et al., 2014) and Transformers (Baly et al., 2020; Liu et al., 2021). But most of these models leverage datasets with ideology labels drawn from a single domain, and it is unclear if any of them can be generalized to diverse genres of text.

Trained on massive quantities of data, Pretrained Language Models (PLMs) have achieved state-of-the-art performance on many text classification problems, with an additional fine-tuning stage on labeled task-specific samples (Devlin et al., 2019; Liu et al., 2019). Though PLMs suggest the promise of

062 generalizable solutions, their ability to acquire the
063 knowledge needed to detect complex features such
064 as ideology from text across genres remains an
065 open question. PLMs have been shown to capture
066 linguistic structures with a *local focus*, such as task-
067 specific words, syntactic agreement, and semantic
068 compositionality (Clark et al., 2019; Jawahar et al.,
069 2019). Although the choice of words is indica-
070 tive of ideology, ideological leaning and stance
071 are often revealed by which entities and events are
072 selected for presentation (Hackett, 1984; Christie
073 and Martin, 2005; Enke, 2020), with the most no-
074 table strand of work in framing theory (Entman,
075 1993, 2007). One such example is demonstrated in
076 Figure 1, where Daily Kos criticizes Trump’s dis-
077 honesty while The Washington Times and Breitbart
078 emphasize the good condition of his health.

079 In this work, we propose to train PLMs for a
080 wide range of ideology-related downstream tasks.
081 We argue that it is critical for PLMs to consider
082 the *global context* of a given article. For instance,
083 as pointed out by Fan et al. (2019), one way to ac-
084 quire such context is through comparison of news
085 articles on the same stories but reported by media
086 of different ideologies. Given the lack of suitable
087 datasets, we first collect a new large-scale dataset,
088 **BIGNEWS**, containing 3,689,229 news articles
089 on politics gathered from 11 US media outlets cov-
090 ering a broad ideological spectrum. We further
091 downsample and then cluster articles in **BIGNEWS**
092 by different media into groups, each group con-
093 sisting of pieces aligned on the same story. The
094 resultant dataset, called **BIGNEWSALIGN**, contains
095 1,060,512 stories with aligned articles.

096 Next we train a new PLM, **POLITICS**¹, based
097 on a Pretraining Objective Leveraging Inter-article
098 Triplet-loss using Ideological Content and Story.
099 Concretely, we leverage continued pretraining (Gu-
100 rurangan et al., 2020), where a novel **ideology ob-**
101 **jective** operating over clusters of articles on the
102 same stories is proposed to compact articles with
103 similar ideology and contrast them with articles
104 of different ideology. The resultant representation
105 can better discern the embedded ideological con-
106 tent. We further enhance it with a **story objective**
107 that ensures the model focuses on meaningful con-
108 tents instead of overly relying on “shortcuts”. Both
109 objectives are used together with our specialized
110 masked language model objective that focuses on
111 entities and sentiments to train **POLITICS**.

¹We will release our data and models upon acceptance.

By experimenting on 11 ideology prediction and
stance detection tasks on 8 datasets of different
genres, including a newly collected dataset from
AllSides, we show that **POLITICS** and its vari-
ants outperform both the strong statistical model-
based comparison SVM and previous PLMs on 10
tasks. Notably, **POLITICS** is especially effective
on long documents, e.g., achieving 10% improve-
ments on both ideology prediction and stance de-
tection tasks over RoBERTa (Liu et al., 2019). This
shows that **POLITICS** can effectively serve as a
general-purpose tool for ideological content analy-
ses. We further show that our model is more robust
in setups with smaller training sets.

2 Related Work

Ideology prediction is one of the core challenges
for understanding political texts, and a critical task
for quantitative political science (Mullins, 1972;
Freeden, 2006; Martin, 2015; Wilkerson and Casas,
2017). Both traditional machine learning meth-
ods (e.g., Naive Bayes, SVM; Evans et al., 2007;
Yu et al., 2008; Sapiro-Gheiler, 2019) and deep
learning models (e.g., RNN; Iyyer et al., 2014)
have been used to predict ideology on a variety of
datasets where ideology labels are available, such
as legislative speeches (Laver et al., 2003) and U.S.
Supreme Court briefs (Evans et al., 2007). Notably,
Liu et al. (2021) pretrains a Transformer-based lan-
guage generator to minimize the ideological bias
in generated text. But generative models are not
as effective as masked language models (MLMs)
at text classification. Therefore, our goal differs in
that we train MLMs that can recognize ideological
content in a wide range of domains and tasks.

Stance Detection. There is a large body of work
on identifying individuals’ stances towards specific
targets from the given text (Thomas et al., 2006;
Walker et al., 2012; Hasan and Ng, 2013). Further-
more, stance detection plays an important role in
measuring public opinions, particularly using eas-
ily accessible posts on social media (Ceron et al.,
2014; Mohammad et al., 2016a; Gautam et al.,
2020; ALDayel and Magdy, 2021). Early stance de-
tection models rely on statistical methods, such as
SVM, based on handcrafted text features (Moham-
mad et al., 2016b; Küçük and Can, 2018). Neu-
ral methods have now been widely investigated
for stance detection, including CNN (Wei et al.,
2016), LSTM (Augenstein et al., 2016), hierarchi-
cal networks (Sun et al., 2018), and unsupervised

representation learning (Darwish et al., 2020).

Recent research focus resides in leveraging PLMs for predicting stances, including incorporating extra features (Prakash and Madabushi, 2020) or distilling knowledge from PLMs (Li et al., 2021). Kawintiranon and Singh (2021) shares a similar spirit with our work by upsampling words for masking. However, they pre-define a list of tokens that are customized for the given targets, which is hard to generalize to new targets. We aim to train PLMs with MLM objectives relying on general-purpose sentiment lexicons and important types of entities, both of which are core elements indicative of stance, to create generalizability of our models.

Pretrained Language Models in the political domain. PLMs, such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and GPT (Radford et al., 2019), have obtained state-of-the-art results on many NLP tasks. With a continued pretraining phase on in-domain data, their predictive performance can be further improved (Gururangan et al., 2020). Based on this idea, domain-specific PLMs, e.g., SciBERT (Beltagy et al., 2019), FinBERT (Yang et al., 2020), LegalBERT (Chalkidis et al., 2020), ClinicalBERT (Huang et al., 2019) and BioBERT (Lee et al., 2020), are trained on curated datasets. However, they all just use the default MLM objective, without considering domain knowledge. In this work, we aim to answer the question: What knowledge needs to be installed into PLMs to produce generalizable tools that can work on various ideology-related tasks? We then design ideology-driven pretraining objectives that allow comparison among articles on the same stories but published by media of different ideologies.

Focusing on the news domain, PLMs have been primarily used for news factuality prediction (Jwa et al., 2019; Zellers et al., 2019; Zhang et al., 2020; Kaliyar et al., 2021) and topic classification (Liu et al., 2020; Büyükköz et al., 2020; Gupta et al., 2020) by fine-tuning on task-specific datasets. We instead target to train new PLMs for usage in broader domains. Furthermore, little has been done for investigating how effectively PLMs can discern political ideology evinced in texts. One exception is Baly et al. (2020), where they also design a triplet-loss pretraining objective to capture ideological content. However, they rely on a smaller dataset consisting of 34,737 articles that are published by the same media but with opposite ideologies, which are scarce. Our pretraining

objective is more practical, and relies on articles aligned as reporting the same stories, but not necessarily from the same outlet. We also release a large-scale dataset, BIGNEWS, for future work in this direction. To the best of our knowledge, we are the first to systematically study and release PLMs for the political domain.

3 Pretraining Datasets

3.1 Data Crawling

We collect pretraining datasets from online news articles with diverse ideological leanings and language usage. We select 11 media outlets based on their ideologies (ranging from far-left to far-right) and popularity.² We then crawl all pages published by them between January 2000 and June 2021, from Common Crawl and Internet Archive. We then follow Raffel et al. (2020) to clean the data, and, additionally, retain news articles related to US politics. Appendix A details cleaning steps for removing non-articles pages, duplicates, non-US politics pages, and boilerplate languages.

The cleaned data, dubbed **BIGNEWS**, contains 3,689,229 political news articles. To mitigate the bias that some media dominate the model training, we downsample the corpus so that each ideology contributes equally. The downsampled corpus, **BIGNEWSBLN**, contains 2,331,552 news articles, with statistics listed in Table 1. We keep 30K as validation. **BIGNEWSBLN** is used to train all baselines and models in this work that employ a MLM objective.

3.2 Aligning Articles on the Same Story

We compare how media outlets from different sides report the same story, which intuitively better captures the ideological content. To this end, we design an algorithm to align articles in **BIGNEWSBLN** that cover the same story. We treat each article as an anchor, and find matches from other outlets based on the following similarity score:

$$\text{sim}(p_i, p_j) = \alpha * \text{sim}_t(p_i, p_j) + (1 - \alpha) * \text{sim}_e(p_i, p_j) \quad (1)$$

where p_i and p_j are two articles, sim_t is the cosine similarity between TF-IDF vectors of p_i and p_j , sim_e is the weighted Jaccard similarity between the sets of named entities³ in p_i and p_j ,

²We use <https://www.allsides.com> and <https://adfontesmedia.com> to decide ideology and <https://www.alexa.com/topsites> to decide popularity.

³Extracted by Stanford CoreNLP (Manning et al., 2014).

	Daily Kos	HPO	CNN	WaPo	NYT	USA Today	AP	The Hill	TWT	FOX	Breitbart
Ideology	L	L	L	L	L	C	C	C	R	R	R
# articles	100,828	241,417	64,988	198,529	173,737	170,737	279,312	322,145	243,181	330,166	206,512
# words	738.7	729.9	655.7	803.2	599.4	691.7	572.3	426.3	522.7	773.5	483.5

Table 1: Statistics of BIGNEWSBLN. Media outlets are sorted by ideology from left (L), center (C), to right (R) based on AllSides and Media Bias Chart. HPO: Huffington Post; WaPo: The Washington Post; NYT: The New York Times; TWT: The Washington Times. Additional statistics of raw size before downsampling, refer to Table A4.

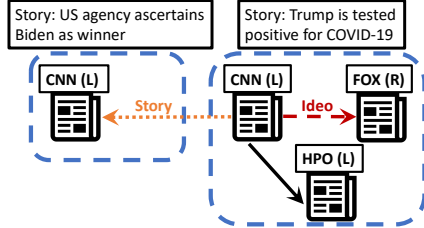


Figure 2: Construction of the ideology and story objectives. The middle CNN article is the anchor in this example. Solid black arrow represents positive-pair relation for both objectives; red dashed arrow denotes negative-pair in ideology objective; orange dashed arrow indicates negative-pair in story objective.

and $\alpha = 0.4$ is a hyperparameter. During alignment, for an article from an outlet to be considered as a match, it must be published within three days before or after the anchor was, has the highest similarity score among articles from the same outlet, and the score is at least $\theta = 0.23$. Hyperparameters α and θ are searched on the Basil dataset (Fan et al., 2019), which contains manually aligned articles from HPO, NYT, and FOX. After deduplicating articles in each story cluster, we get **BIGNEWSALIGN**, containing 1,060,512 clusters with an average of 4.29 articles in each. Appendix B details the alignment algorithm.

4 POLITICS with Continued Pretraining

Here we introduce our continued pretraining methods based on a newly proposed **ideology objective** that drives representation learning to better discern ideological content by comparing same-story articles (§4.1), which is further augmented by a **story objective** to better focus on the content. They are combined with the masked language model objective which is tailored to focus on entities and sentiments (§4.2) to produce **POLITICS** (§4.3).

4.1 Ideology-driven Pretraining Objectives

To promote representation learning that better captures ideological content, we leverage **BIGNEWSALIGN** with articles grouped by sto-

ries to provide story-level background for model training. Concretely, we use triplet loss that operates over **triplets** of $\langle \text{anchor}, \text{positive}, \text{negative} \rangle$ (Schroff et al., 2015), to encourage anchor and positive samples to have closer representations compared to those of anchor and negative samples.

Our primary pretraining objective, i.e., ideology objective, uses the triplet loss to teach the model to acquire **ideology-informed representations** by comparing same-story articles written by media of different ideologies. As shown in Figure 2, given a story cluster, we choose an article published by media on the left or right as the *anchor*. We then take articles in the cluster with the same ideology as *positive* samples, and articles with the opposite ideology as *negative* ones. The ideology objective is formulated as follows:

$$\mathcal{L}_{\text{ideo}} = \sum_{t \in \mathcal{T}_{\text{ideo}}} \left[\left\| \mathbf{t}^{(a)} - \mathbf{t}^{(p)} \right\|_2 - \left\| \mathbf{t}^{(a)} - \mathbf{t}^{(n)} \right\|_2 + \delta_{\text{ideo}} \right]_+ \quad (2)$$

where $\mathcal{T}_{\text{ideo}}$ is the set of all possible ideology triplets in the training set, $\mathbf{t}^{(a)}$, $\mathbf{t}^{(p)}$, and $\mathbf{t}^{(n)}$ are the [CLS] representations of anchor, positive, and negative articles in triplet t , δ_{ideo} is a hyperparameter, and $[\cdot]_+$ is defined as $\max(\cdot, 0)$.

Next, we augment the ideology objective with a story objective to allow the model to focus on **semantically meaningful content** and to prevent the model from focusing on “shortcuts” (such as media-specific languages) to detect ideology. To construct story triplets, we use the same $\langle \text{anchor}, \text{positive} \rangle$ pairs as in the ideology triplet, and then treat articles from the same media outlet but on different stories as negative samples. Similarly, our story objective is formulated as follows:

$$\mathcal{L}_{\text{story}} = \sum_{t \in \mathcal{T}_{\text{story}}} \left[\left\| \mathbf{t}^{(a)} - \mathbf{t}^{(p)} \right\|_2 - \left\| \mathbf{t}^{(a)} - \mathbf{t}^{(n)} \right\|_2 + \delta_{\text{story}} \right]_+ \quad (3)$$

where $\mathcal{T}_{\text{story}}$ contains all story triplets in training, and δ_{story} is a hyperparameter searched on the validation set.

4.2 Entity- and Sentiment-aware MLM

Here we present a specialized MLM objective to collaborate with our triplet loss based objectives for better representation learning. Notably, political framing effect is often reflected in which entities are selected for reporting (Gentzkow et al., 2019). Moreover, the occurrence of sentimental content along with the entities also signal stances (Mohammad et al., 2016b). Therefore, we take a masking strategy that upsamples *entity* tokens (Sun et al., 2019; Guu et al., 2020; Kawintiranon and Singh, 2021) and *sentiment* words to be masked for the MLM objective, which improves from prior pre-training work that only considers article-level comparison (Baly et al., 2020).

Concretely, we consider named entities of PERSON, NORP, ORG, GPE and EVENT types. We detect sentiment words using lexicons by Hu and Liu (2004) and Wilson et al. (2005). To focus MLM training more on entities and sentiment, we mask them with a 30% probability, and then randomly mask remaining tokens until 15% (the same probability as used in BERT) of all tokens are reached. We also follow BERT on replacing masked tokens with [MASK], random, and original tokens.

4.3 Overall Pretraining Objective

We combine the aforementioned objectives as our final pretraining objective as follows:

$$\mathcal{L} = \beta * \mathcal{L}_{\text{ideology}} + \gamma * \mathcal{L}_{\text{story}} + (1 - \beta - \gamma) * \mathcal{L}_{\text{MLM}} \quad (4)$$

where $\beta = 0.25$ and $\gamma = 0.25$.

Using \mathcal{L} , POLITICS is produced via continued training on RoBERTa-base⁴ (Liu et al., 2019). Details of hyperparameters are listed in Table A5.

5 Experiments

Given the importance of ideology prediction and stance detection tasks in political science (Thomas et al., 2006; Wilkerson and Casas, 2017; Chatsiou and Mikhaylov, 2020), we conduct extensive experiments on a wide spectrum of datasets with 11 tasks (§5.1). We then compare with baselines of both traditional machine learning models and prior PLMs (§5.2), and among our model variants (§5.3). We present and discuss results in §5.5, where POLITICS outperform all baselines on 8 out of 11 tasks.

All models that use a MLM objective are pre-trained with BIGNEWSBLN, and the ones with

⁴We use roberta-base model card from Huggingface.

Data	Genre	# Train	Len.	Split
Congress Speech (Gentzkow et al., 2018)	speech	7,000	538	rand.
AllSides (newly collected)	news	7,878	863	time
BASIL-article (Fan et al., 2019)	news	450	693	story
BASIL-sentence (Fan et al., 2019)	news	1,197	27	story
Hyperpartisan (Kiesel et al., 2019)	news	425	556	rand.
VAST (Allaway and McKeown, 2020)	cmt	11,545	102	rand.*
YouTube User (Wu and Resnick, 2021)	cmt	1,114	1,213	user
YouTube Cmt (Wu and Resnick, 2021)	cmt	6,832	197	user
SemEval (Mohammad et al., 2016a)	tweet	2,251	17	rand.*
Twitter (Preoțiu-Pietro et al., 2017)	tweet	1,079	2,298	user

Table 2: Datasets used for evaluating PLMs vary in text genre, training set size (# Train), length, and split criterion. Time split means training on the “past” data and test on the “future”. *: splits by the original work.

our ideology and story objectives are pretrained on BIGNEWSALIGN.

5.1 Datasets and Tasks

Our tasks are discussed below, with dataset statistics listed in Table 2. Please refer to Appendix D for dataset processing details.

Ideology prediction tasks are evaluated on the following datasets.

- Congress Speech (**CongS**; Gentzkow et al., 2018) contains parsed speeches from US congressional records, each labeled as liberal or conservative.
- AllSides⁵ (**AIS**) is a website that assesses political bias and ideology of US media outlets. In this study, we collect articles from AllSides with their ideological leanings on a 5-point scale.
- Hyperpartisan (**HP**; Kiesel et al., 2019) is a shared task of identifying news that takes an extreme left-wing or right-wing standpoint.
- YouTube (Wu and Resnick, 2021) contains discussions on YouTube. **YT (cmt.)** and **YT (user)** refer to predicting left or right at comment- and user-level (by concatenating comments by the same user).
- Twitter (**TW**; Preoțiu-Pietro et al., 2017) collects a group of Twitter users with self-reported ideologies on a 7-point scale. Each user is represented by their posted tweets.

Stance detection tasks, that predict a subject’s attitude (positive, negative, neutral) towards a given target from a piece of text, are listed below.

- BASIL (Fan et al., 2019) contains news articles with annotations on authors’ stances towards given entities. **BASIL (sent.)** and **BASIL (art.)**

⁵<https://www.allsides.com>.

	Ideology Prediction						Stance Detection						All avg	
	YT (cmt.)	CongS	HP	AllS	YT (user)	TW	Ideo. avg	SEval (seen)	SEval (unseen)	Basil (sent.)	VAST	Basil (art.)		Stan. avg
SVM	65.34	<u>71.31</u>	61.25	52.51	66.49	42.85	59.96	51.18	32.89	51.08	39.54	30.77	41.09	51.38
BERT	64.64	<u>65.88</u>	48.42	60.88	65.24	44.20	58.21	65.07	40.39	62.81	70.53	45.61	56.88	57.61
RoBERTa	66.72	67.25	60.43	74.75	67.98	48.90	64.34	70.15	63.08	68.16	76.25	41.36	63.80	64.09
Our models with triplet loss objective only														
Ideology Obj.	66.20	<u>68.18</u>	<u>64.15</u>	76.52	<u>68.15</u>	42.66	64.31	68.78	59.61	64.18	76.03	<u>44.94</u>	62.71	63.58
Story Obj.	66.09	<u>69.11</u>	56.70	74.59	<u>68.89</u>	46.53	63.65	69.02	63.54	67.21	<u>76.66</u>	53.16	<u>65.92</u>	<u>64.68</u>
Ideology Obj. + Story Obj.	<u>68.91</u>	<u>69.10</u>	<u>63.08</u>	<u>76.23</u>	<u>77.58</u>	48.98	<u>67.31</u>	69.66	<u>63.17</u>	64.37	76.18	<u>47.01</u>	<u>64.08</u>	<u>65.84</u>
Our models with masked language model objective only														
Random	67.82	70.32	60.59	73.54	70.77	44.62	64.61	69.16	60.39	69.94	77.11	39.16	63.15	63.95
Upsamp. Ent.	69.06	<u>70.32</u>	60.09	70.89	<u>71.40</u>	<u>47.16</u>	<u>64.82</u>	<u>69.81</u>	<u>63.08</u>	69.49	76.76	<u>46.46</u>	<u>65.12</u>	<u>64.96</u>
Upsamp. Sentiment	67.41	70.03	56.05	72.35	<u>74.93</u>	<u>48.15</u>	<u>64.82</u>	<u>70.09</u>	<u>60.81</u>	<u>71.28</u>	76.61	<u>44.42</u>	<u>64.64</u>	<u>64.74</u>
Upsamp. Ent. + Sentiment	<u>68.31</u>	71.42	58.02	71.90	<u>71.04</u>	<u>47.31</u>	<u>64.67</u>	<u>69.25</u>	<u>62.84</u>	69.23	77.10	<u>43.16</u>	<u>64.32</u>	<u>64.51</u>
POLITICS	<u>67.83*</u>	70.86	70.25*	<u>74.93</u>	78.73*	<u>48.92</u>	68.59*	69.41	61.26	73.41*	<u>76.73</u>	<u>51.94*</u>	66.55*	67.66*

Table 3: Macro F1 scores on 11 evaluation tasks (average of 5 runs). Tasks are sorted by text length, short to long, within each group. “All avg” is the average of all 11 tasks. **Best** results are in bold and second best results are underlined. Our models with triplet-loss objectives that outperform RoBERTa are in blue. Our models with specialized sampling methods that outperform with vanilla MLM (Random) are in green. POLITICS uses both Ideology Obj. + Story Obj. and Upsamp. Ent. + Sentiment. Results where POLITICS outperforms all baselines are highlighted in red, with * indicating statistical significance (t -test, $p \leq 0.05$).

are prediction tasks on sentence and article-level.

- VAST (Allaway and McKeown, 2020) collects online comments from The New York Times “Room for Debate” section, with stances labeled towards the debate topic.
- SemEval (Mohammad et al., 2016a) is a shared task on detecting stances in tweets. We consider two setups to predict on seen, i.e. SEval (seen), and unseen, i.e., SEval (unseen), entities.

5.2 Baselines

We consider three baselines. First, we train a linear SVM using unigram and bigram features for each task, since it is a common baseline in political science (Yu et al., 2008; Diermeier et al., 2012b). Hyperparameters and feature selection are described in Table A7. We further compare with BERT and RoBERTa, following the standard fine-tuning process for ideology prediction tasks and using the prompt described in §5.4 for stance detection.

5.3 Our Model Variants

We consider several variants for POLITICS. First, using triplet loss objective only, we report results by our models that are trained on BIGNEWSALIGN with ideology objective (*Ideology Obj.*), story objective (*Story Obj.*), and both.

Next, we continue pretaining RoBERTa with MLM objective only, using vanilla MLM objective (*Random*), entity focused objective (*Upsamp. Ent.*), sentiment focused objective (*Upsamp. Sentiment*), as well as upsampling both entity and sentiment.

5.4 Fine-tuning Procedure

We fine-tune each neural model for up to 10 epochs, with early stopping enabled. We select the best fine-tuned model on validation sets using F1. Details of experimental setups are in Table A6.

Ideology Prediction. We follow common practice of using the [CLS] token for standard fine-tuning (Devlin et al., 2019). For Twitter and YouTube User data, we encode them using a sliding window and aggregate by mean pooling.

Stance Detection. We follow Schick and Schütze (2021) on using prompts to fine-tune models for stance detection. We curate 11 prompts (in Table A8) and choose the best one based on the average F1 by RoBERTa on all stance detection tasks, as shown below:

p [SEP] *The stance towards {target} is* [MASK] .

The model is trained to predict [MASK] for stance, conditioned on the input p and {target}.

5.5 Main Results

Table 3 presents the average F1 scores on all tasks. POLITICS achieves the best overall average F1 score across the board, 3.6% better than the strongest baseline, RoBERTa. More importantly, POLITICS alone outperforms all the baselines on 8 out of 11 tasks, including more than 10% of improvement for ideology labeling on Hyperpartisan and Youtube user-level. We attribute the performance gain to our proposed



Figure 3: Macro F1 aggregated over tasks of different formality, training size, document length and aggregation method (single post vs. user posts). Results show that POLITICS performs better on handling formal language, small training sets, and longer text.

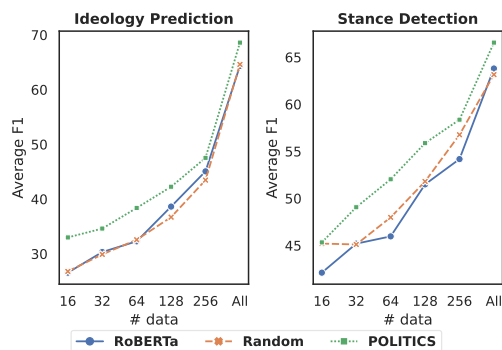


Figure 4: Average of ideology prediction and stance detection performances with few-shot learning. POLITICS uniformly outperforms RoBERTa that is continued pretrained with vanilla MLM (Random).

ideology-driven pretraining objective which helps capture partisan content.

Moreover, our ideology-driven objectives helps acquire knowledge needed to discern ideology as well as stance detection. When equipping the RoBERTa model with ideology and story objectives but no MLM objective, it achieves the second best overall performance on ideology prediction and also improves on stance detection tasks.

Next, focusing on entities better identifies stance. Simply continuing training RoBERTa with vanilla MLM objective (Random) does not lead to performance gain on stance detection. However, leveraging upsampling methods makes a difference. By increasing the sampling ratio of entities, stance-related tasks are improved by 2%.

On Texts of Different Characteristics. Based on Table 2, we can further study the model’s performance based on different properties of the data: language formality, size of the training set, document length, and aggregation level. As shown in Figure 3, using each property (with definition listed in Appendix E), we divide the tasks into two categories. POLITICS yields greater improvements on more formal and longer text, since pretraining is done on news articles. POLITICS is also more robust to training sets with small sizes, showing the potential effectiveness in few-shot learning, which is echoed by our study in §6.1.

6 Further Analyses

6.1 Few-shot Learning

We first fine-tune all PLMs on small numbers of samples. As shown in Figure 4, we find that POLITICS performs consistently better than the two counterparts on both tasks when the training

set is small. More importantly, naively training RoBERTa on the large BIGNEWSBLN does not help ideology prediction. By contrast, our ideology-driven objective helps to better learn ideology, e.g., when using only 16 samples for fine-tuning on the ideology tasks, compared to the baselines.

6.2 Ablation Study on POLITICS

We show the impact of removing each ideology-driven pretraining objective and upsampling strategy from POLITICS in Table 4. First, removing the ideology objective results in the most loss on both tasks, again, demonstrating the effectiveness of our triplet-loss formulation over same-story articles. Removing the story objective also hurts the overall performance by 1% but improves the ideology prediction marginally. This shows that the story objective functions as an auxiliary constraint to avoid over-fitting on the “shortcuts” for discerning ideologies. Moreover, removing upsampling strategies generally weakens POLITICS’s performance, but only to a limited extent.

We also experiment with a setup with hard-ideology learning (i.e., directly predicting the ideology of each article without using triplet-loss objectives). Not surprisingly, this variant (POLITICS +Ideo. Pred.) outperforms POLITICS on ideology prediction since it can directly learn ideology from the annotated labels. However, it has been overfitted to ideology prediction tasks and loses the generalizability of transferring knowledge, thus yields the worse performance on stance detection.

6.3 Visualizing Attentions

On the Hyperpartisan task, we visualize the last layer’s attention weights between the [CLS] token and all other tokens by POLITICS

	Ideology Prediction							Stance Detection						All avg
	YT (cmt.)	CongS	HP	AllS	YT (user)	TW	Ideo. avg	SEval (seen)	SEval (unseen)	Basil (sent.)	VAST	Basil (art.)	Stan. avg	
POLITICS	67.83	70.86	70.25	74.93	78.73	48.92	68.59	69.41	61.26	73.41	76.73	51.94	66.55	67.66
No Ideology Obj.	-3.78	-2.17	-16.35	-3.28	-12.54	-3.43	-6.93	-0.38	-0.83	-4.22	-0.45	-16.01	-4.38	-5.77
No Story Obj.	+1.98	+0.64	-0.72	+0.70	+0.29	-1.78	+0.19	-1.23	+2.94	-3.36	-0.87	-10.75	-2.66	-1.11
No Upsamp. Ent.	+0.18	-0.65	-0.05	+0.55	-0.29	-1.20	-0.24	+0.62	-0.67	-3.74	-0.55	-1.20	-1.11	-0.64
No Upsamp. Sentiment	+0.75	-0.28	+0.22	-1.27	-0.11	-1.40	-0.35	-0.84	+1.67	-3.91	-1.10	+1.44	-0.55	-0.44
POLITICS + Ideo. Pred.	+1.46	+1.10	-1.01	+4.72	+2.02	-3.96	+0.72	+0.41	-0.52	-3.82	+0.12	-3.10	-1.38	-0.23

Table 4: Ablation study results on POLITICS. POLITICS + Ideo. Pred.: triplet-loss objective is replaced with a hard label prediction objective on ideology of articles (left vs. right). **Best** results are in bold. Darker red shows greater improvements. Darker blue indicates larger performance drop. The ideology objective contributes the most to POLITICS, followed by the story objective.

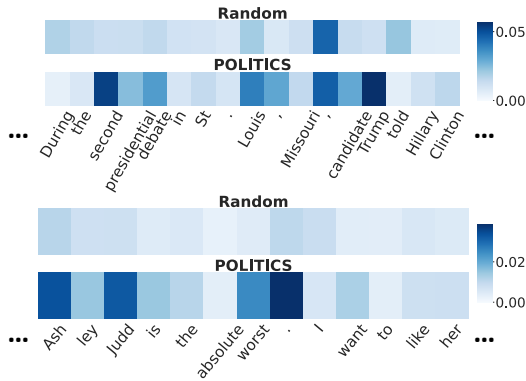


Figure 5: Last layer attention scores between [CLS] token and other input tokens (aggregated over all heads). In the first sentence, POLITICS captures “presidential debate” and “Trump”. In the second sentence, POLITICS captures “worst” and “Ashley Judd”. Longer versions of the plots are in Figures A1 and A2.

and RoBERTa pretrained with vanilla MLM on BIGNEWSBLN. We observe that POLITICS is able to capture salient entities and sentiments in the text, such as “Trump”, “Ashley Judd”, “presidential debate”, and “the worst”, as illustrated in Figure 5. This finding confirms that our ideology-driven objective and upsampling strategies can help the model focus more on entities of political interest as well as better recognize sentiments. More examples can be found in Appendix F.

6.4 POLITICS on Different Ideologies

Finally, we measure whether PLMs would acquire ideological bias as measured by whether they fit with languages used by a specific ideology. Concretely, we evaluate PLMs on 30K held-out articles of different ideologies from BIGNEWSBLN with perplexity. As illustrated in Figure 6, while MLM objective (*Random*) is effective at fitting a corpus, i.e., having the lowest perplexities, we observe that triplet-loss objectives acts as a regularization during pretraining, shown by the similar perplexities

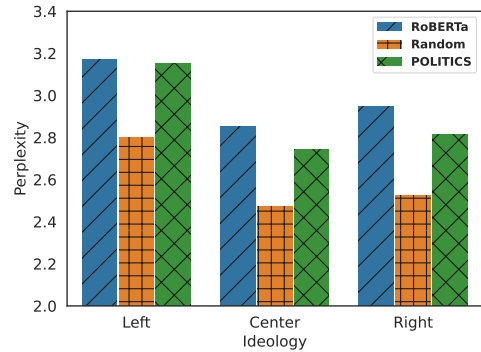


Figure 6: Perplexities of different models on 30K validation articles of different ideologies in BIGNEWSBLN. Perplexities do not drop on POLITICS, suggesting it can yield superior predictive performance while not overfitting with ideological languages.

between POLITICS and RoBERTa. Interestingly, we find center and right articles have lower perplexity than that of left articles. We hypothesize that it relates to the findings in political science that during the recent period of political polarization in the US. Republicans have become somewhat more coherent and similar than Democrats (Grossmann and Hopkins, 2016; Benkler et al., 2018), making them easier to predict.

7 Conclusion

We study the problem of training general-purpose tools for ideology content understanding and prediction. We present POLITICS, trained with novel ideology-driven pretraining objectives based on the comparisons of same-story articles written by media outlets of different ideologies. To facilitate the model training, we also collect a large-scale dataset, BIGNEWS, consisting of news articles of different ideological leanings. Experiments on diverse datasets for the tasks of ideology prediction and stance detection show that POLITICS outperforms strong baselines, even with a limited amount of labeled samples for training.

8 Ethical Considerations

8.1 BIGNEWS Collection

All news articles were collected in a manner consistent with the terms of use of the original sources and the intellectual property and the privacy rights of the original authors of the texts (i.e., source owners). In the data collection process, the collectors honored privacy rights of article authors and no sensitive information was collected (e.g., writers’ identifications). All participants involved in the data collection process have completed human subjects research training at their affiliated institutions. We also consulted Section 107⁶ of the U.S. Copyright Act and ensured that our collection action fell under fair use category.

8.2 Dataset Usage

BIGNEWS will be released under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.⁷ Pretraining corpus details are included in Section 3. The other eight datasets used for downstream evaluation are obtained in the following two ways: **CongS**, **HP**, **BASIL**, **VAST** and **SEval** are obtained by direct download. For **YT** and **TW**, we consult with the corresponding authors and then obtain the datasets from them with verbal agreement on not sharing the dataset to the public. We further crawl **AIS** data from AllSides website while complying with terms of use. Dataset details are listed in Section 5.1 and Appendix D.

8.3 Benefit and Potential Misuse of BIGNEWS and POLITICS

Intended use. Assisting the general public to automatically measure ideology of diverse genres of texts. For example, POLITICS can help the general public know where their representatives stand on key issues. Our experiments in Section 5 matches how POLITICS would be deployed in real life when handling both ideology prediction and stance detection. We believe that our extensive experiments have covered the major usage of POLITICS.

Failure mode is defined as a situation where POLITICS fails to correctly predict the ideology of an individual or an input text. Ideally, the interpretation of our model’s prediction should be carried out

within the broader context of the input text. However, when taken out of context, prediction results may be misinterpreted by users.

Potential harms. No known harms are observed if POLITICS is being used as intended and functioning correctly. However, if POLITICS malfunctions on stance detection tasks, it could generate opposite results, which might deliver misinformation or make users misunderstand a political figure’s stance towards a policy. For vulnerable populations (e.g., people who cannot make the right judgements), the harm might be tremendously magnified if they fail to interpret the ideology prediction and stance detection results in an expected way or blindly trust machine responses.

Misuse potential. Users may mistakenly take the machine prediction as a golden rule or a fact. We would recommend any politics-related machine learning models put up an “use with caution” message to encourage users to check more resources or consult political science experts to reduce the risk of being misled by one single source.

Bias Mitigation. In our data preprocessing step, we downsample BIGNEWS to BIGNEWSBLN to ensure that each ideology contributes equally to the corpus that is later used for continued pretraining, with the purpose of minimizing potential bias. We do not think that POLITICS explicitly encodes any bias. In Figure 6, the discrepancy in perplexities among different ideologies is more related to the greater coherence among Republicans than Democrats (Grossmann and Hopkins, 2016; Benkler et al., 2018), rather than POLITICS encoding biased knowledge.

In conclusion, there is no greater than minimal risk/harm introduced by either BIGNEWSBLN or POLITICS. However, to discourage the misuse, we will always warn users that model predictions are for informational purpose only and users should always resort to the broader context to reduce the risk of absorbing biased information.

⁶<https://www.copyright.gov/title17/92chap1.html#107>

⁷<https://creativecommons.org/licenses/by-nc-sa/4.0/>

662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717

References

Abeer ALDayel and Walid Magdy. 2021. Stance detection on social media: State of the art and trends. *Information Processing & Management*, 58(4):102597.

Emily Allaway and Kathleen McKeown. 2020. Zero-Shot Stance Detection: A Dataset and Model using Generalized Topic Representations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8913–8931, Online. Association for Computational Linguistics.

Stephen Ansolabehere, Jonathan Rodden, and James M Snyder. 2008. The strength of issues: Using multiple measures to gauge preference stability, ideological constraint, and issue voting. *American Political Science Review*, 102(2):215–232.

Isabelle Augenstein, Tim Rocktäschel, Andreas Vlachos, and Kalina Bontcheva. 2016. Stance detection with bidirectional conditional encoding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 876–885. The Association for Computational Linguistics.

Ramy Baly, Giovanni Da San Martino, James R. Glass, and Preslav Nakov. 2020. We can detect your bias: Predicting the political ideology of news articles. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 4982–4991. Association for Computational Linguistics.

Pablo Barberá, John T. Jost, Jonathan Nagler, Joshua A. Tucker, and Richard Bonneau. 2015. Tweeting from left to right: Is online political communication more than an echo chamber? *Psychological Science*, 26(10):1531–1542.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3613–3618. Association for Computational Linguistics.

Yochai Benkler, Robert Faris, and Hal Roberts. 2018. *Network propaganda: Manipulation, disinformation, and radicalization in American politics*. Oxford University Press.

Adam Bonica. 2013. Mapping the ideological marketplace. *ERN: Models of Political Processes: Rent-Seeking*.

Berfu Büyüköz, Ali Hürriyetoğlu, and Arzucan Özgür. 2020. Analyzing ELMo and DistilBERT on socio-political news classification. In *Proceedings of the Workshop on Automated Extraction of Socio-political Events from News 2020*, Marseille, France. European Language Resources Association (ELRA).

Shuyang Cao and Lu Wang. 2021. Inference time style control for summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Online. Association for Computational Linguistics. 718
719
720
721
722
723

Andrea Ceron, Luigi Curini, Stefano M Iacus, and Giuseppe Porro. 2014. Every tweet counts? how sentiment analysis of social media can improve our knowledge of citizens’ political preferences with an application to italy and france. *New Media & Society*, 16(2):340–358. 724
725
726
727
728
729

Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: the muppets straight out of law school. *CoRR*, abs/2010.02559. 730
731
732
733

Kakia Chatsiou and Slava Jankin Mikhaylov. 2020. Deep learning for political science. *CoRR*, abs/2005.06540. 734
735
736

Frances Christie and James R Martin. 2005. *Genre and institutions: Social processes in the workplace and school*. A&C Black. 737
738
739

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of bert’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@ACL 2019, Florence, Italy, August 1, 2019*, pages 276–286. Association for Computational Linguistics. 740
741
742
743
744
745
746
747

Kareem Darwish, Peter Stefanov, Michaël J. Aupetit, and Preslav Nakov. 2020. Unsupervised user stance detection on twitter. In *Proceedings of the Fourteenth International AAAI Conference on Web and Social Media, ICWSM 2020, Held Virtually, Original Venue: Atlanta, Georgia, USA, June 8-11, 2020*, pages 141–152. AAAI Press. 748
749
750
751
752
753
754

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics. 755
756
757
758
759
760
761
762
763
764

Daniel Diermeier, Jean-François Godbout, Bei Yu, and Stefan Kaufmann. 2012a. Language and ideology in congress. *British Journal of Political Science*, 42(1):31–55. 765
766
767
768

Daniel Diermeier, Jean-François Godbout, Bei Yu, and Stefan Kaufmann. 2012b. Language and ideology in congress. *British Journal of Political Science*, 42(1):31–55. 769
770
771
772

773	Benjamin Enke. 2020. What you see is all there is. <i>The Quarterly Journal of Economics</i> , 135(3):1363–1398.	826
774		827
775	Robert M. Entman. 1993. Framing: Toward clarification of a fractured paradigm. <i>Journal of Communication</i> , 43(4):51–58.	828
776		829
777		830
778	Robert M. Entman. 2007. Framing bias: Media in the distribution of power. <i>Journal of Communication</i> , 57(1):163–173.	831
779		832
780		833
781	Michael Evans, Wayne McIntosh, Jimmy Lin, and Cynthia Cates. 2007. Recounting the courts? applying automated content analysis to enhance empirical legal research. <i>Journal of Empirical Legal Studies</i> , 4(4):1007–1039.	834
782		835
783		836
784		837
785		838
786	Lisa Fan, Marshall White, Eva Sharma, Ruisi Su, Prafulla Kumar Choubey, Ruihong Huang, and Lu Wang. 2019. In plain sight: Media bias through the lens of factual reporting. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 6342–6348. Association for Computational Linguistics.	839
787		840
788		841
789		842
790		843
791		844
792		845
793		846
794		847
795		848
796	Morris P. Fiorina and Samuel J. Abrams. 2008. Political polarization in the american public. <i>Annual Review of Political Science</i> , 11(1):563–588.	849
797		850
798		851
799	Michael Freeden. 2006. Ideology and political theory. <i>Journal of Political Ideologies</i> , 11(1):3–22.	852
800		853
801	Matthew J Gabel and John D Huber. 2000. Putting parties in their place: Inferring party left-right ideological positions from party manifestos data. <i>American Journal of Political Science</i> , pages 94–103.	854
802		855
803		856
804		857
805	Akash Gautam, Puneet Mathur, Rakesh Gosangi, Debanjan Mahata, Ramit Sawhney, and Rajiv Ratn Shah. 2020. # metooma: Multi-aspect annotations of tweets related to the metoo movement. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 14, pages 209–216.	858
806		859
807		860
808		861
809		862
810		863
811	Matthew Gentzkow, Jesse M. Shapiro, and Matt Taddy. 2018. Congressional record for the 43rd-114th congresses: Parsed speeches and phrase counts.	864
812		865
813		866
814	Matthew Gentzkow, Jesse M. Shapiro, and Matt Taddy. 2019. Measuring group differences in high-dimensional choices: Method and application to congressional speech. <i>Econometrica</i> , 87(4):1307–1340.	867
815		868
816		869
817		870
818	Tim Groseclose, Steven D Levitt, and James M Snyder. 1999. Comparing interest group scores across time and chambers: Adjusted ada scores for the us congress. <i>American political science review</i> , 93(1):33–50.	871
819		872
820		873
821		874
822		875
823	Matt Grossmann and David A Hopkins. 2016. <i>Asymmetric politics: Ideological Republicans and group interest Democrats</i> . Oxford University Press.	876
824		877
825		878
	Shloak Gupta, S Bolden, Jay Kachhadia, A Korsunskaya, and J Stromer-Galley. 2020. Polibert: Classifying political social media messages with bert. In <i>Social, Cultural and Behavioral Modeling (SBP-BRIMS 2020) conference. Washington, DC</i> .	879
		880
	Suchin Gururangan, Ana Marasovic, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don’t stop pretraining: Adapt language models to domains and tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 8342–8360. Association for Computational Linguistics.	879
		880
	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Realm: Retrieval-augmented language model pre-training.	879
		840
	Robert A Hackett. 1984. Decline of a paradigm? bias and objectivity in news media studies. <i>Critical Studies in Media Communication</i> , 1(3):229–259.	841
		842
	Kazi Saidul Hasan and Vincent Ng. 2013. Stance classification of ideological debates: Data, models, features, and constraints. In <i>Sixth International Joint Conference on Natural Language Processing, IJCNLP 2013, Nagoya, Japan, October 14-18, 2013</i> , pages 1348–1356. Asian Federation of Natural Language Processing / ACL.	843
		844
	Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In <i>Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’04</i> , page 168–177, New York, NY, USA. Association for Computing Machinery.	845
		846
	Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. 2019. Clinicalbert: Modeling clinical notes and predicting hospital readmission. <i>CoRR</i> , abs/1904.05342.	847
		848
	Mohit Iyyer, Peter Enns, Jordan L. Boyd-Graber, and Philip Resnik. 2014. Political ideology detection using recursive neural networks. In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers</i> , pages 1113–1122. The Association for Computer Linguistics.	849
		850
	Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. What does BERT learn about the structure of language? In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 3651–3657. Association for Computational Linguistics.	851
		852
	Heejung Jwa, Dong Bin Oh, Kinam Park, Jang Kang, and Hueiseok Lim. 2019. exbake: Automatic fake news detection model based on bidirectional encoder representations from transformers (bert). <i>Applied Sciences</i> .	853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880

881	Rohit Kumar Kaliyar, Anurag Goswami, and Pratik Narang. 2021. Fakebert: Fake news detection in social media with a bert-based deep learning approach. <i>Multim. Tools Appl.</i> , 80(8):11765–11788.	937
882		938
883		939
884		940
885	Kornraphop Kawintiranon and Lisa Singh. 2021. Knowledge enhanced masked language model for stance detection. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4725–4735, Online. Association for Computational Linguistics.	941
886		942
887		943
888		944
889		945
890		946
891		947
892	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. 2019. SemEval-2019 task 4: Hyperpartisan news detection. In <i>Proceedings of the 13th International Workshop on Semantic Evaluation</i> , pages 829–839, Minneapolis, Minnesota, USA. Association for Computational Linguistics.	948
893		949
894		950
895		951
896		952
897		953
898		954
899		955
900	Heemin Kim and Richard C Fording. 1998. Voter ideology in western democracies, 1946–1989. <i>European Journal of Political Research</i> , 33(1):73–97.	956
901		957
902		958
903	Dilek Küçük and Fazli Can. 2018. Stance detection on tweets: An svm-based approach. <i>CoRR</i> , abs/1803.08910.	959
904		960
905		961
906	Michael Laver, Kenneth Benott, and John Garry. 2003. Extracting policy positions from political texts using words as data. <i>American Political Science Review</i> , 97(2):311–331.	962
907		963
908		964
909		965
910	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. <i>Bioinform.</i> , 36(4):1234–1240.	966
911		967
912		968
913		969
914		970
915	Jeffrey Lewis, Keith Poole, Howard Rosenthal, Adam Boche, Aaron Rudkin, and Luke Sonnet. 2021. Voteview: Congressional roll-call votes database.	971
916		972
917		973
918	Yang Li, Yuqing Sun, and Nana Zhu. 2021. Berttoenn: Similarity-preserving enhanced knowledge distillation for stance detection. <i>Plos one</i> , 16(9):e0257130.	974
919		975
920		976
921	Jingang Liu, Chunhe Xia, Xiaojian Li, Haihua Yan, and Tengting Liu. 2020. A bert-based ensemble model for chinese news topic prediction. In <i>Proceedings of the 2020 2nd International Conference on Big Data Engineering</i> , New York, NY, USA. Association for Computing Machinery.	977
922		978
923		979
924		980
925		981
926		982
927	Ruibo Liu, Chenyan Jia, Jason Wei, Guangxuan Xu, Lili Wang, and Soroush Vosoughi. 2021. Mitigating political bias in language models through reinforced calibration. In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 14857–14866. AAAI Press.	983
928		984
929		985
930		986
931		987
932		988
933		989
934		990
935		991
936		992
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized bert pretraining approach.	993
		994
	Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In <i>Association for Computational Linguistics (ACL) System Demonstrations</i> , pages 55–60.	995
		996
	John Levi Martin. 2015. What is ideology? <i>Sociologia, Problemas e Práticas</i> , pages 9–31.	997
		998
	Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiao-Dan Zhu, and Colin Cherry. 2016a. Semeval-2016 task 6: Detecting stance in tweets. In <i>Proceedings of the 10th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2016, San Diego, CA, USA, June 16-17, 2016</i> , pages 31–41. The Association for Computer Linguistics.	999
		1000
	Saif M. Mohammad, Parinaz Sobhani, and Svetlana Kiritchenko. 2016b. Stance and sentiment in tweets. <i>CoRR</i> , abs/1605.01655.	1001
		1002
	Willard A. Mullins. 1972. On the concept of ideology in political science. <i>American Political Science Review</i> , 66(2):498–510.	1003
		1004
	Andrew Peterson and Arthur Spirling. 2018. Classification accuracy as a substantive quantity of interest: Measuring polarization in westminster systems. <i>Political Analysis</i> , 26(1):120–128.	1005
		1006
	Keith T Poole and Howard Rosenthal. 1985. A spatial model for legislative roll call analysis. <i>American journal of political science</i> , pages 357–384.	1007
		1008
	Anushka Prakash and Harish Tayyar Madabushi. 2020. Incorporating count-based features into pre-trained models for improved stance detection. <i>CoRR</i> , abs/2010.09078.	1009
		1010
	Daniel Preotiu-Pietro, Ye Liu, Daniel Hopkins, and Lyle Ungar. 2017. Beyond binary labels: Political ideology prediction of Twitter users. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 729–740, Vancouver, Canada. Association for Computational Linguistics.	1011
		1012
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	1013
		1014
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer.	1015
		1016

990	Eitan Sapiro-Gheiler. 2019. Examining political trustworthiness through text-based measures of ideology. In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 10029–10030. AAAI Press.	1046
991		1047
992		1048
993		1049
994		1050
995		1051
996		1052
997		1053
998		
999	Timo Schick and Hinrich Schütze. 2021. Exploiting cloze questions for few shot text classification and natural language inference.	1054
1000		1055
1001		1056
1002		1057
1003	Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. <i>2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	1058
1004		1059
1005		1060
1006		1061
1007	Boris Shor and Nolan McCarty. 2011. The ideological mapping of american legislatures. <i>American Political Science Review</i> , 105(3):530–551.	1062
1008		1063
1009		1064
1010	Valentin I. Spitzkovsky and Angel X. Chang. 2012. A cross-lingual dictionary for English Wikipedia concepts. In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)</i> , Istanbul, Turkey. European Language Resources Association (ELRA).	1065
1011		1066
1012		1067
1013		
1014		
1015		
1016	Qingying Sun, Zhongqing Wang, Qiaoming Zhu, and Guodong Zhou. 2018. Stance detection with hierarchical attention network. In <i>Proceedings of the 27th International Conference on Computational Linguistics, COLING 2018, Santa Fe, New Mexico, USA, August 20-26, 2018</i> , pages 2399–2409. Association for Computational Linguistics.	1068
1017		1069
1018		1070
1019		
1020		
1021		
1022		
1023	Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. Ernie: Enhanced representation through knowledge integration.	1071
1024		1072
1025		1073
1026		1074
1027	Matt Thomas, Bo Pang, and Lillian Lee. 2006. Get out the vote: Determining support or opposition from congressional floor-debate transcripts. In <i>EMNLP 2006, Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing, 22-23 July 2006, Sydney, Australia</i> , pages 327–335. ACL.	1075
1028		1076
1029		1077
1030		1078
1031		1079
1032		1080
1033	Andrea Volkens, Tobias Burst, Werner Krause, Pola Lehmann, Theres MatthieÄÿ, Nicolas Merz, Sven Regel, Bernhard WeÄÿels, and Lisa Zehnter. 2021. The manifesto data collection. manifesto project (mrg/cmp/marpor). version 2021a.	1081
1034		1082
1035		
1036		
1037		
1038	Marilyn A. Walker, Pranav Anand, Rob Abbott, and Ricky Grant. 2012. Stance classification using dialogic properties of persuasion. In <i>Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 3-8, 2012, Montréal, Canada</i> , pages 592–596. The Association for Computational Linguistics.	1083
1039		1084
1040		1085
1041		1086
1042		1087
1043		1088
1044		
1045		
	Wan Wei, Xiao Zhang, Xuqin Liu, Wei Chen, and Tengjiao Wang. 2016. pkudblab at SemEval-2016 task 6 : A specific convolutional neural network system for effective stance detection. In <i>Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)</i> , pages 384–388, San Diego, California. Association for Computational Linguistics.	1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
	John Wilkerson and Andreu Casas. 2017. Large-scale computerized text analysis in political science: Opportunities and challenges. <i>Annual Review of Political Science</i> , 20(1):529–544.	1054
		1055
		1056
		1057
	Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2005. Recognizing contextual polarity in phrase-level sentiment analysis. In <i>Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing</i> , pages 347–354, Vancouver, British Columbia, Canada. Association for Computational Linguistics.	1058
		1059
		1060
		1061
		1062
		1063
		1064
	Siqi Wu and Paul Resnick. 2021. Cross-partisan discussions on youtube: Conservatives talk to liberals but liberals don't talk to conservatives.	1065
		1066
		1067
	Yi Yang, Mark Christopher Siy Uy, and Allen Huang. 2020. Finbert: A pretrained language model for financial communications. <i>CoRR</i> , abs/2006.08097.	1068
		1069
		1070
	Bei Yu, Stefan Kaufmann, and Daniel Diermeier. 2008. Classifying party affiliation from political speech. <i>Journal of Information Technology & Politics</i> , 5(1):33–48.	1071
		1072
		1073
		1074
	Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 9051–9062.	1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
	Tonglu Zhang, Di Wang, Huanhuan Chen, Zhiwei Zeng, Wei Guo, Chunyan Miao, and Lizhen Cui. 2020. Bdann: Bert-based domain adaptation neural network for multi-modal fake news detection. <i>2020 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–8.	1083
		1084
		1085
		1086
		1087
		1088

Filter Patterns	
url	/video/, /gallery/, /slideshow/
title	weekly digest, 10 sites you should know, day’s end roundup, photos of the week, 5 things you need to know

Table A1: Examples of patterns used to filter out pages that are not news articles.

Appendix A BIGNEWS Cleaning Steps

In this section, we provide the details of our data cleaning steps for BIGNEWS. We adopt the following cleaning steps to only keep news articles that relate to US politics in BIGNEWS.

Remove Non-article Pages. Online news websites post many contents that are not news articles. We remove such pages by checking the page title and url. We create a list of patterns to filter out invalid pages. Some examples of the patterns are shown in Table A1.

Remove Duplicate Pages. We use the character level edit distance to find duplicate pages. Specifically, we use the following formula to calculate the difference between page a and page b :

$$\text{diff}(a, b) = \text{dist}(a, b) / \max(\text{len}(a), \text{len}(b)) \quad (5)$$

where $\text{dist}(a, b)$ is the Levenshtein distance between a and b . If the difference is less than 0.1, we consider two pages as duplicates of each other. For duplicated pages, we only keep the one that has the earliest publish date. Following this procedure, we remove duplicated pages in each media outlet.

Remove Non-politics Pages. To filter out non-politics pages, we build a politics classifier to check whether a page is about politics or not. We create the training data from BIGNEWS. Because the url is usually indicative of the content of a page, we use keywords in the url to retrieve politics and non-politics training data. The lists of keywords are shown in Table A2. This results in a training dataset with 400,462 politics pages and 310,377 non-politics pages. We also randomly sample 888 pages from the remaining dataset and manually annotate them to use as the test set.

With the training data, we train a unigram and bigram TF-IDF vectorizer to extract features and a Logistic Regression model to do classification. Because the lists of keywords in Table A2 might not

Keywords	
Politics	/politics/, /political/, /policy/, /election/, /elections/, /allpolitics/
Non-politics	/travel/, /sports/, /life/, /movie/, /entertainment/, /science/, /music/, /plated/, /leisure/, /showbiz/, /lifestyle/, /fashion/, /art/

Table A2: Keywords used to retrieve positive and negative training data.

Url Keywords	Text US Keywords
/world/, /international/, /europe/, /africa/, /asia/, /latin-america/, /middle-east/	U.S., United States, Obama, Trump, Bush, Biden, Pompeo, Clinton, Pence

Table A3: Examples of keywords used to filter out non-US pages. For text keywords, we include all presidents, vice presidents, and secretaries of state of the US after 2000.

be complete, we use the trained classifier to classify remaining pages and add those that are classified with high confidence⁸ to the training data. This results in a larger training set with 957,424 politics pages and 987,898 non-politics pages. We train the final classifier on the larger training set and achieve a 88.67% F-1 score and 88.18% accuracy on the test data.

Remove Non-US Pages. We filter out pages that do not relate to the US by looking for keywords in the url. We create a list of keywords that identify potential non-US pages. For those pages, we further check if they contain US related keywords and only remove those that have no US related keywords. Examples of keywords we use are shown in Table A3.

Remove Media-info Leaking Phrases. To prevent the model from learning features specific to individual media outlets, we adopt two cleaning steps. First, we mask phrases that mention the media outlet (e.g., New York Times, NYTimes, and nytimes.com). Second, we create a list of patterns for frequently appearing sentences (more than 100 times) of each media outlet. For example, for the following sentence: “author currently serves as a senior political analyst for [MASK] Channel and

⁸We use 0.95 for politics pages and 0.1 for non-politics pages.

	# article before downsample	Earliest article	Latest article
Daily Kos	235,244	2009-01-02	2021-06-30
HuffPost (HPO)	560,581	2000-11-30	2021-06-30
CNN	152,579	2000-01-01	2021-06-30
The Washington Post (WaPo)	461,032	2000-01-01	2021-06-30
The New York Times (NYT)	403,191	2000-01-01	2021-06-22
USA Today	174,525	2001-01-01	2021-06-30
Associated Press (AP)	285,685	2000-01-01	2021-06-30
The Hill (Hill)	337,256	2002-10-06	2021-06-30
The Washington Times (TWT)	336,056	2000-01-01	2021-06-30
Fox News (FOX)	457,550	2001-01-12	2021-06-25
Breitbart News (Breitbart)	285,530	2009-01-08	2021-06-30

Table A4: Statistics of BIGNEWS corpus. Media outlets are sorted by ideology from left to right.

contributes to all major political coverage.” Both the author name and the sentence itself can leak media outlet information. As such sentences usually appear at the beginning or end of the article, we remove the first and last two paragraphs that contain any of the patterns.

Appendix B Story Alignment

As shown in Equation 1, we combine text similarity and entity similarity as the final similarity score. We only consider the title and the first five sentences in the calculation. We further require aligned articles a and b to satisfy two constraints:

- Difference in publish dates of a and b is less than or equal to three days.
- a and b must contain at least one common named entity in the title and first three sentences.

We use CoreNLP to extract named entities in the article (Manning et al., 2014). For constraint two, we further use Crosswikis to map the entity to a unique concept in Wikipedia (Spitkovsky and Chang, 2012). When calculating entity similarity, we split the entity into single words and remove stop words. After alignment, we use the procedure described in Appendix A to remove duplicate articles in the same story cluster. The final hyperparameters we use are $\alpha = 0.4$ and $\theta = 0.23$.

Evaluate Alignment Algorithm To evaluate the performance of the alignment algorithm, we use the AllSides dataset collected in Cao and Wang (2021). The dataset consists of manually aligned news articles from 251 media outlets. After removing media outlets not in the BigNews corpus, we have 2,904 articles on 1,316 events. We add the

AllSides dataset into the BigNews corpus and use each AllSides article as the anchor article for the alignment algorithm. We use the aligned article in the AllSides dataset as the relevant article and the algorithm achieves 0.679 mean reciprocal rank.

Appendix C Continued Pretraining and Fine-tuning

C.1 Continued Pretraining

We initialize all variants of POLITICS with RoBERTa-base model (Liu et al., 2019), which contains about 125M parameters. We train each model using 8 Quadro RTX 8000 GPUs for 2,500 steps. The total training time for POLITICS is 20 hours. For other variants of POLITICS, the training time could be shorter. Table A5 lists out the training hyperparameters.

Training Details. For triplet loss objectives, we only consider triplets in each mini-batch. We skip a batch if it contains no triplet. For MLM objective, we truncate the article if it has more than 512 tokens. When masking entities and sentiment words, we only consider those with at most five tokens.

C.2 Fine-tuning

For both ideology prediction and stance detection tasks, we fine-tune each model for up to 10 epochs. We use early stopping and select the best checkpoint on validation set among 10 epochs. For ideology prediction tasks, we follow standard practice of using [CLS] token and feedforward neural networks (FNN) for classification. For stance detection tasks, we use prompts to fine-tune PLMs. We curate 11 prompts as shown in Table A8. We

Hyperparameter	Value
number of steps	2,500
batch size	2048
maximum learning rate	0.0005
learning rate scheduler	linear decay with warmup
warmup percentage	6%
optimizer	AdamW
weight decay	0.01
AdamW beta weights	0.9, 0.98
δ_{ideo}	0.5
δ_{story}	1.0

Table A5: Hyperparameters used in continued pretraining.

select the best prompt based on the performance of RoBERTa. Fine-tuning hyperparameters are listed in Table A6. Hyperparameters of SVM classifier are listed in Table A7.

Appendix D Downstream Evaluation Datasets

This section lists more details of the eight datasets used in our downstream evaluation as well as their processing steps.

D.1 Ideology Prediction

- Congress Speech (**CongS**; Gentzkow et al., 2018): We filter out speeches with less than 80 words, and we use the party affiliation of the speaker as the ideology of the speech.
- AllSides⁹ (**AIS**): We crawl articles from AllSides and use the annotated ideology of media outlets as the ideology of articles. We further annotate ideology of each article by the ideology of the media outlet.
- Hyperpartisan (**HP**; Kiesel et al., 2019): We convert the benchmark into a 3-way classification task by projecting media-level ideology annotations to article level.
- YouTube (Wu and Resnick, 2021) contains cross-partisan discussions between liberals and conservatives on YouTube. In our experiments, we only keep controversial comments: 1) A video must have at least 1,500 comments and 150,000 views; 2) A comment must have at least 20 replies. The original dataset annotates users'

⁹<https://www.allsides.com>.

Hyperparameter	Value
number of epochs	10
patience	4
maximum learning rate	0.00001 or 0.00002
learning rate scheduler	linear decay with warmup
warmup percentage	6%
optimizer	AdamW
weight decay	0.001
AdamW beta weights	0.9, 0.999
# FNN layer	2
hidden layer dimension in FNN	768
dropout in FNN	0.1
sliding window size	512
sliding window overlap	64

Table A6: Hyperparameters used to fine-tune PLMs.

Hyperparameter	Value
kernel	linear
regularization strength	0.3, 1, or 3
features	unigram and bigram TF-IDF
minimum document frequency	5
maximum document frequency	$0.7 * D $

Table A7: Hyperparameters used to train SVM. $|D|$ is the number of documents in the training set.

- ideology on a 7-point scale. We further convert it into a 3-way classification task that contains left, center, and right ideologies. For comment-level prediction task on **YT (cmt.)**, we annotate the ideology by the user-level ideology which is provided. For user-level prediction on **YT (user)**, we concatenate all comments from a user.
- Twitter (**TW**; PreoŃiuc-Pietro et al., 2017): We crawl recent tweets for each user and remove replies and non-English tweets. We assume users' ideologies do not change after their self-report since prior work has shown that people's ideology is less likely to change across the political spectrum (Fiorina and Abrams, 2008). We sort all tweets from a user by time and concatenate them.

D.2 Stance Detection

- **BASIL** (Fan et al., 2019): We convert the original dataset such that the new tasks are to predict the stance towards a target at two granularities:

Prompt	Verbalizer
<i>p</i> [SEP] <i>The stance towards {target} is</i> [MASK].	negative or positive
<i>p</i> [SEP] <i>It reveals a</i> [MASK] <i>stance on</i> {target}.	negative or positive
<i>p</i> [SEP] <i>The speaker holds a</i> [MASK] <i>attitude towards</i> {target}.	negative or positive
<i>p</i> [SEP] <i>What is the stance on</i> {target}? [MASK].	Negative or Positive
<i>p</i> [SEP] <i>The previous passage</i> [MASK] {target}.	opposes or favors
<i>p</i> [SEP] <i>The stance on</i> {target} <i>is</i> [MASK].	negative or positive
<i>p</i> [SEP] <i>The stance towards</i> {target}: [MASK].	negative or positive
<i>p</i> [SEP] <i>The author</i> [MASK] {target}.	opposes or favors
<i>p</i> [SEP] [MASK] {target}	oppose or favor
<i>p</i> [SEP] [MASK] . {target}	No or Yes
<i>p</i> [SEP] [MASK] {target}	No or Yes

Table A8: List of prompts designed for stance detection tasks. *p* is the input text. {target} is the target of interests. Verbalizer maps the label (against) to the token (negative) that we want models to predict. Some datasets have a third label (neutral).

article (**art.**) and sentence (**sent.**) levels. The targets in the dataset can be a person (e.g., Donald Trump) or an organization (e.g., Justice Department).

- VAST (Allaway and McKeown, 2020) is the task to predict the stance of a comment towards a target. The targets in the dataset are noun phrases covering a broad range of topics (e.g., immigration and home schoolers).
- SemEval (SEval; Mohammad et al., 2016a) is the task to predict a tweet’s stance towards a target where a target at test time could be seen or unseen during training. The dataset contains six targets: Atheism, Climate Change, Feminist, Hillary Clinton, Abortion, and Donald Trump (unseen).

Appendix E Task Property

This section introduces detailed definitions of four properties, i.e., how we divide tasks into two categories for each property.

- Formality: Speech and news genres are considered formal while the remainder are informal.
- Training set size: Datasets with more than 2,000 training samples are considered large, otherwise small.
- Document length: Datasets with average document length larger than 500 are considered “long” while the remainder are short.
- Aggregation level: If a dataset is a collection of single articles/posts/tweets, then it is categorized into “Single”; If posts are concatenated and aggregated at the user-level, then it is marked as “User”. Specifically, only YouTube User and

Twitter in Table 2 fit into “User” category.

Appendix F Visualize Attention Weights

In this section, we visualize attention weights for more examples.

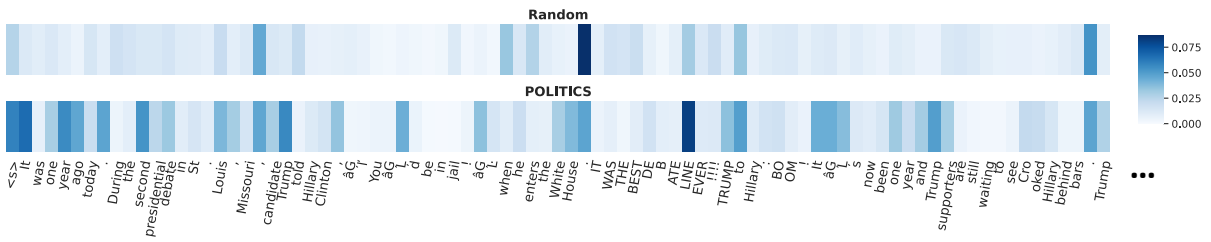


Figure A1: Example 1. Last layer attention weights between [CLS] token and other tokens in the input. We illustrate the first 85 tokens of the article.

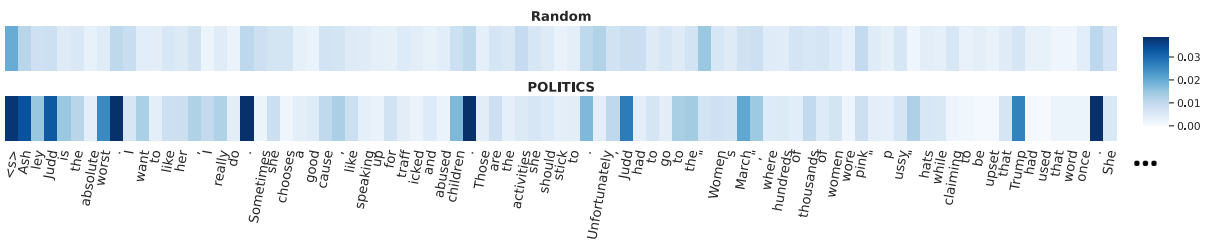


Figure A2: Example 2. Last layer attention weights between [CLS] token and other tokens in the input. We illustrate the first 85 tokens of the article.

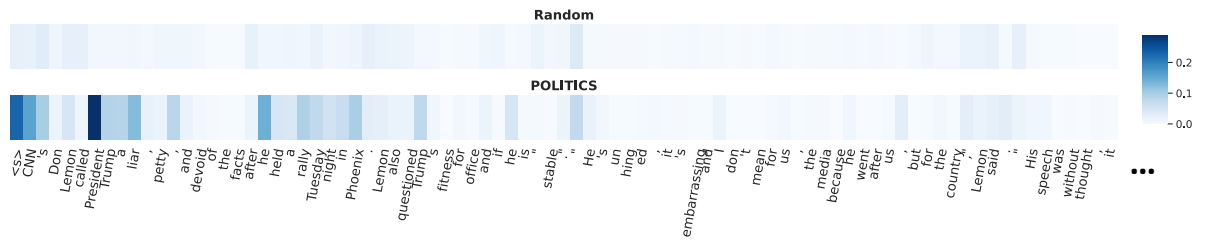


Figure A3: Example 3. Last layer attention weights between [CLS] token and other tokens in the input. We illustrate the first 85 tokens of the article.

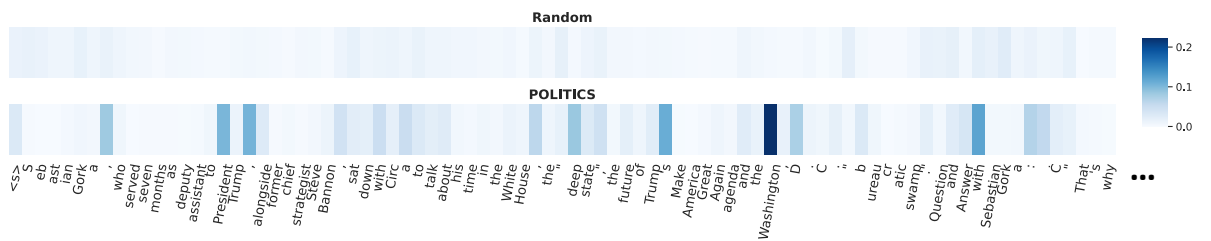


Figure A4: Example 4. Last layer attention weights between [CLS] token and other tokens in the input. We illustrate the first 85 tokens of the article.

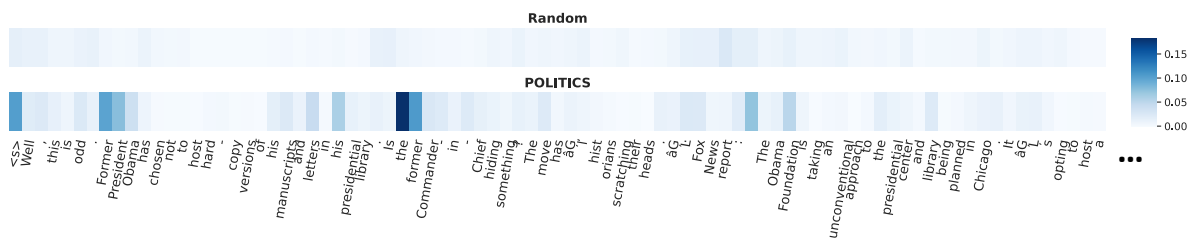


Figure A5: Example 5. Last layer attention weights between [CLS] token and other tokens in the input. We illustrate the first 85 tokens of the article.