

# HQP: A Human-Annotated Dataset for Detecting Online Propaganda

Abdurahman Maarouf\* Dominik Bär\* Dominique Geissler\* Stefan Feuerriegel\*

\*Munich Center for Machine Learning (MCML) & LMU Munich  
[{a.maarouf,baer,d.geissler,feuerriegel}@lmu.de}](mailto:{a.maarouf,baer,d.geissler,feuerriegel}@lmu.de)

## Abstract

Online propaganda poses a severe threat to the integrity of societies. However, existing datasets for detecting online propaganda have a key limitation: they were annotated using *weak* labels that can be noisy and even incorrect. To address this limitation, our work makes the following contributions: (1) We present HQP: a novel dataset ( $N = 30,000$ ) for detecting online propaganda with *high-quality* labels. To the best of our knowledge, HQP is the first dataset for detecting online propaganda that was created through human annotation. (2) We show empirically that state-of-the-art language models fail in detecting online propaganda when trained with weak labels (AUC: 64.03). In contrast, state-of-the-art language models can accurately detect online propaganda when trained with our high-quality labels (AUC: 92.25), which is an improvement of  $\sim 44\%$ . (3) To address the cost of labeling, we extend our work to few-shot learning. Specifically, we show that prompt-based learning using a small sample of high-quality labels can still achieve a reasonable performance (AUC: 80.27). Finally, we discuss implications for the NLP community to balance the cost and quality of labeling. Crucially, our work highlights the importance of high-quality labels for sensitive NLP tasks such as propaganda detection.

**Disclaimer:** Our work contains potentially offensive language and manipulative content. Reader’s discretion is advised.

## 1 Introduction

Propaganda refers to communication that is primarily used to influence, persuade, or manipulate public opinions (Smith, 2022). Nowadays, propaganda is widely used as practices of modern warfare (e.g., in the ongoing Russo-Ukrainian war) and thus poses a significant threat to the integrity of societies (Kowalski, 2022). In this regard, social

media presents especially fertile grounds for disseminating propaganda at an unprecedented scale.

Existing NLP works for propaganda detection generally focus on different types of content. One literature stream aims at detecting propaganda in official news (e.g., Rashkin et al., 2017; Barrón-Cedeño et al., 2019; Da San Martino et al., 2019) but fails to detect propaganda in online content (Wang et al., 2020). Another literature stream aims at detecting propaganda in online content from social media (Wang et al., 2020; Vijayaraghavan and Vosoughi, 2022). Here, existing works rely upon datasets that were exclusively annotated using *weak* labels and were, therefore, not validated by humans. Because of this, labels can be noisy and even incorrect. We later provide empirical support for this claim and show that the overlap between weak labels and human annotations is only  $\sim 41\%$ . Given that social media is increasingly used for propaganda dissemination (Geissler et al., 2022), a tailored dataset is needed to tackle it. To the best of our knowledge, there is *no* dataset for detecting online propaganda that was constructed through human annotation and validation and thus comprises *high-quality* labels.

To fill this gap, we develop HQP: a novel dataset with **high-quality** labels for detecting online propaganda. HQP consists of  $N = 30,000$  tweets in English from the Russo-Ukrainian war. We use human annotation and validation to generate *high-quality* labels. Specifically, we follow best practices and make use of a rigorous multi-annotator, multi-batch procedure for annotation (Song et al., 2020). We then leverage state-of-the-art, pre-trained language models (PLMs), i.e., BERT, RoBERTa, and BERTweet, to benchmark the performance in detecting online propaganda using weak labels vs. our high-quality labels. We find that high-quality labels are crucial for performance. We further acknowledge that human annotation also incurs labeling costs, and, to address

this, we extend our work to few-shot learning (i.e., prompt-based learning).

Our main **contributions** are as follows:<sup>1</sup>

1. We construct HQP, a novel dataset with high-quality labels for online propaganda detection using human-annotated labels.
2. We show that PLMs for detecting online propaganda using high-quality labels outperform PLMs using weak labels by a large margin.
3. We adapt few-shot learning for online propaganda detection by prompting PLMs.

## 2 Related Work

**Detecting harmful content:** Prior literature in NLP has aimed to detect a broad spectrum of harmful content such as hate speech (e.g., Badjatiya et al., 2017; Mathew et al., 2021; Pavlopoulos et al., 2022), rumors (e.g., Zhou et al., 2019; Bian et al., 2020; Xia et al., 2020; Wei et al., 2021), and fake news (e.g., Zellers et al., 2019; Liu et al., 2020; Lu and Li, 2020; Jin et al., 2022). Further, claim detection has been studied, for example, in the context of the Russo-Ukrainian war (La Gatta et al., 2023). Overall, literature for detecting harmful content makes widespread use of datasets that were created through human annotations (e.g., Founta et al., 2018; Thorne et al., 2018; Mathew et al., 2021), yet outside of online propaganda. We add by constructing a dataset through human annotations that is tailored to online propaganda.

**Detecting propaganda content:** Previous works for propaganda detection can be loosely grouped by the underlying content, namely (1) official news and (2) social media. We briefly review both in the following.

(1) *News*. To detect propaganda in official news, existing works leverage datasets that originate from propagandistic and non-propagandistic news outlets (Rashkin et al., 2017; Barrón-Cedeño et al., 2019; Da San Martino et al., 2019; Solopova et al., 2023), yet these datasets are not tailored to online content from social media. As a case in point, Wang et al. (2020) previously examined the capability of machine learning to transfer propaganda detection between news and online content, yet found challenges in doing so.

(2) *Social media*. To detect propaganda in social media, existing works create datasets from

online platforms such as Twitter. For example, TWE (Wang et al., 2020) combines a random sample of tweets (representing the non-propagandistic class) with a sample of tweets from the Internet Research Agency (representing the propagandistic class). TWEETSPIN (Vijayaraghavan and Vosoughi, 2022) is a dataset with tweets that are annotated with weak labels along different types of propaganda techniques by mining accusations in the replies and quotes to each tweet. Notably, all existing datasets for detecting online propaganda were created through weak annotation (see Table 1). To this end, labels can oftentimes be noisy or even incorrect. We fill this void by developing a human-annotated dataset for online propaganda detection.

Methodologically, earlier works generally rely upon feature engineering (tf-idf) (Barrón-Cedeño et al., 2019) and LSTMs (Rashkin et al., 2017; Wang et al., 2020). More recently, PLMs such as BERT evolved as the state-of-the-art method to detect propaganda (Da San Martino et al., 2019; Vijayaraghavan and Vosoughi, 2022). Later, we thus also adopt state-of-the-art PLMs to assess the role of weak vs. high-quality labels.

**Few-shot learning in NLP:** Generally, constructing large-scale datasets with high-quality labels in NLP is costly. Hence, there is a growing interest in few-shot learning. Common methods typically leverage prompting, where the downstream task is reformulated to resemble the masked language modeling task the PLM was trained on (e.g., Radford et al., 2019; Brown et al., 2020; Gao et al., 2021; Schick and Schütze, 2021; Liu et al., 2023). Prompting has been highly successful in few-shot learning, e.g., for rumor detection (Lin et al., 2023) and humor detection (Li et al., 2023). However, to the best of our knowledge, no work has so far adapted few-shot learning to detect propaganda.

## 3 Dataset Construction (HQP)

In the following, we construct a human-annotated dataset of English social media content with propaganda (HQP). For this, we construct a corpus of tweets with Russian propaganda from the 2022 invasion of Ukraine. We collect tweets from February 2021 until October 2022, i.e., our timeframe starts one year before the invasion due to the widespread opinion that the invasion was planned

<sup>1</sup>Code and data publicly available at <https://github.com/abdumaa/HiQualProp>.

Dataset	Domain	Level	Human ann.	Model	Few-shot
Rashkin et al. (2017)	News	Document	✗	LSTM	✗
Barrón-Cedeño et al. (2019)	News	Document	✗	Maximum entropy classifier	✗
Da San Martino et al. (2019)	News	Fragment	✓	Multi-granularity network	✗
Solopova et al. (2023)	News	Document	✗	BERT	✗
Wang et al. (2020) ("TWE")	Social media	Short-text	✗	LSTM	✗
Vijayaraghavan and Vosoughi (2022) ("TWEETSPIN")	Social media	Short-text	✗	Multi-view transformer	✗
<b>HQP (ours)</b>	Social media	Short-text	✓	BERT, RoBERTa, BERTweet	✓

Table 1: Overview of existing datasets for propaganda detection aimed at (i) news and (ii) online content.

far in advance.<sup>2</sup> We intentionally choose the Russo-Ukrainian war due to its significance for world politics (Kowalski, 2022) and the sheer size of the propaganda campaign (Geissler et al., 2022). Our methodology for constructing HQP follows best practices for human annotation (Song et al., 2020). Specifically, we build upon a four-step process: (1) data collection, (2) sampling, and (3) human annotation.

### 3.1 Data Collection

In practice, social media content with propaganda is rare in comparison to non-propaganda (e.g., well below 0.1%). Therefore, simply collecting a random subset of tweets will contain only very few samples from the positive class. Instead, we follow the methodology in Founta et al. (2018) and perform a stratified search. Thereby, we separately generate candidates for the positive class ( $D_+$ ) and for the negative class ( $D_-$ ) as shown in Figure 1. The latter are context-related samples in that they also discuss topics related to the Russo-Ukrainian war but are likely not propaganda. Thereby, we eventually capture a challenging setting in which we can evaluate how accurately propaganda and non-propaganda can be discriminated.

**(i) Candidate search for positive class ( $D_+$ ):** Analogous to Vijayaraghavan and Vosoughi (2022), we expect that propaganda on Twitter is often called out in replies or quotes (e.g., some users debunk propaganda as such). We thus access the Twitter Historical API and perform a keyword-based search. Specifically, we crawl replies and quotes that contain phrases (keywords) that may accuse the original tweet of propaganda, such as ("russian"  $\wedge$  "propaganda") or ("war"  $\wedge$  "propaganda"). The full list is in Table A.1 in the supplements.

We create the list of search keywords through

<sup>2</sup><https://www.nytimes.com/2021/04/09/world/europe/russia-ukraine-war-troops-intervention.html>

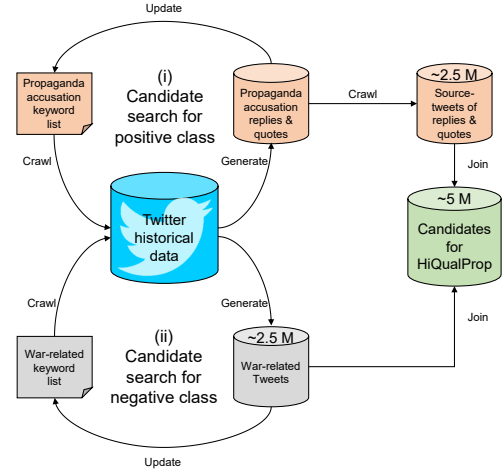


Figure 1: Data collection of candidate tweets for HQP.

an iterative process: (1) In each iteration, the current list of keywords is used to filter for English-language replies and quotes from Twitter. (2) We then manually scan the most frequent words (including bi- and tri-grams) for phrases that can potentially qualify as propaganda accusations. (3) We add these to our list of keywords. We repeat the process for three iterations and use the final list of search keywords to retrieve our set of replies and quotes. Afterward, we crawl the corresponding source tweets, which resulted in  $\sim 2.5$  million candidates for  $D_+$ .

**(ii) Candidates for negative class ( $D_-$ ):** To collect candidates for the negative class, we crawl a random sample of 2.5 million tweets that discuss the Russo-Ukrainian war but that have not necessarily been identified as propaganda through users. For this, we use a similar iterative procedure to generate a keyword list as for the positive class. However, we now perform a keyword search only for source tweets (but not for replies or quotes). Example keywords are ("russia"  $\wedge$  "war") and ("ukraine"  $\wedge$  "war"). The complete list of search keywords is in Table A.2 in the supplements.

**Postprocessing:** We postprocess the candidates

for both the positive class ( $D_+$ ) and the negative class ( $D_-$ ). Specifically, we filter out duplicates, non-English tweets, and very short tweets (i.e., fewer than 5 words). The resulting union of both postprocessed candidate sets contains  $|D_+| + |D_-| \sim 3.2$  million samples.

Unlike Vijayaraghavan and Vosoughi (2022), we do not perform weak labeling by simply assigning a label to a tweet depending on whether it is in  $D_+$  or  $D_-$ , respectively. Instead, we generate high-quality labels through human annotation. This is motivated by our observation that many samples in  $D_+$  cover the Russo-Ukrainian war but do not qualify as propaganda. Hence, weak labeling would lead to many false positives.

### 3.2 Boosted Sampling

We collect  $N = 30,000$  tweets from the postprocessed union of  $D_+$  and  $D_-$  for human annotation. We adopt boosted sampling (Founta et al., 2018) as we observe that the majority of samples from the previous step cover ‘normal’ content and do not necessarily qualify as propaganda. Our research objective requires that there is a sufficient proportion of positive class labels in our dataset, since, otherwise, the dataset will not be useful for the research community. To address the class imbalance, we use weighted sampling. For this, we generate weights  $w_i$  for each tweet as the inverse term frequency of potential propaganda-related phrases (e.g., “nobody talks about”), that is,

$$w_i = \frac{n_i}{\sum_{j=1}^M n_j}, \quad (1)$$

where  $n_i$  is the number of unique occurrences of propaganda-related phrases and  $M = |D_+| + |D_-|$ . Here, we use the list of 189 potential propaganda phrases from Vijayaraghavan and Vosoughi (2022). As a result, our boosted sampling approach will increase the likelihood that actual propaganda content (true positives) is later annotated (rather than false positives).

### 3.3 Human Annotation

To annotate our data, we recruit human workers from Prolific (<https://www.prolific.co/>) to label the tweets as propagandistic or non-propagandistic. Workers are pre-selected according to strict criteria: residency in the UK/US, English as a first language; enrollment in an undergraduate, graduate, or doctoral degree; a minimum approval rate of 95%; and min. 500 completed submissions

on Prolific. The annotation instructions are in Appendix B of the supplement.

Workers are asked to annotate two labels for each tweet. The first is a binary label ( $BL$ ) to classify propagandistic vs. non-propagandistic content. The second is a propaganda-strategy label ( $PSL$ ) aimed at capturing the context-related strategy behind propagandistic tweets. Therefore, if a tweet is annotated as propagandistic, the worker is asked to select one of four context-related propaganda strategies that are used in this tweet (thus giving  $PSL$ ). The four propaganda strategies were carefully chosen after manually studying a sample of 2,000 tweets and discussing different options with an expert team of propaganda researchers. Specifically, for  $PSL$ , workers have to decide whether the propagandistic tweet is designed to influence opinions (1) against Western countries, (2) against Ukraine, (3) pro-Russian government, or (4) aimed at other countries.

Our annotation follows a batch procedure according to best practices (Song et al., 2020), i.e., a pool of workers annotates a subset of the data to avoid fatigue. We thus split the dataset ( $N = 30,000$ ) into 300 batches with 100 tweets each. Each batch is annotated by two workers. Beforehand, we manually annotate 10 tweets of each batch with respect to  $BL$  to measure the quality of the annotations: If (a) a worker incorrectly labeled more than 20% of the internally annotated tweets or (b) the inter-annotator agreement between both workers has a Cohen’s kappa (Cohen, 1960)  $\leq 0.4$ , we discard the annotation and repeat the annotation for the batch. Overall, we had to discard and redo 7.5% of the batch annotations. When annotators disagreed on  $BL$  for individual tweets, we resolve the conflicts as follows: the  $BL$  is then re-annotated by randomly assigning it to one of the top 25 annotators. If there is disagreement on the  $PSL$  after resolving the disagreement on the  $BL$ , the final  $PSL$  is decided by the author team. The latter was the case for only 2.6% of the tweets. Altogether, this corroborates the reliability of our multi-annotator, multi-batch procedure for annotation.

We further crawled additional meta information about authors, which we report in Appendix F.

### 3.4 Dataset Statistics

Table 2 reports summary statistics for our HQP. Table 3 lists five example tweets and their corresponding labels, i.e.,  $BL$  and  $PSL$ .



	Propaganda = true	Propaganda = false	Overall
Num. of tweets ( $N$ )	4,610	25,390	30,000
Avg. tweet length (in chars)	238.71	216.90	220.25
Num. of unique authors	3,910	20,140	23,317

Table 2: Summary statistics for HQP.

## 4 Methods

In our experiments, we follow state-of-the-art methods from the literature (see Sec. 2) to ensure the comparability of our results. To this end, we use a binary classification task (propaganda = 1, otherwise = 0). Results are reported as the average performance over five separate runs. In each run, we divide the dataset into train (70%), val (10%), and test (20%) using a stratified shuffle split. All experiments are conducted on a Ubuntu 20.04 system, with 2.30 GHz Intel Xeon Silver 4316 CPU and two NVIDIA A100-PCIE-40GB GPUs.

### 4.1 Fine-Tuning PLMs

#### 4.1.1 PLMs

We use the following PLMs in our experiments: **BERT-large** (Devlin et al., 2019), **RoBERTa-large** (Liu et al., 2019), and **BERTweet-large** (Nguyen et al., 2020). The latter uses the pretraining procedure from RoBERTa but is tailored to English tweets in order to better handle social media content. We report implementation details in Appendix C.<sup>3</sup>

#### 4.1.2 Baselines

We compare the fine-tuning procedure on our high-quality labels vs. baselines that make use of weak labels. All evaluations are based on separate test splits of HQP with human verification.

- **TWE:** We fine-tune on weak labels of the public TWE dataset (Wang et al., 2020).

- **TWEETSPIN:** The TWEETSPIN dataset with weak labels (Vijayaraghavan and Vosoughi, 2022) is not public, and we thus replicate the data collection procedure ( $N = 3,223,867$ ).

- **HQP-weak:** As an ablation study, we construct a dataset with weak labels based on our data collection procedure from HQP. Specifically, we map from our classification into  $D_+$  and  $D_-$ . We generate  $N = |\text{HQP}|$  samples so that the size is comparable to HQP. This allows us later to isolate the effect of the dataset size from the role of weak vs

high-quality labels.<sup>4</sup>

### 4.2 Prompt-Based Learning

For few-shot learning, we leverage state-of-the-art prompt-based learning (Liu et al., 2023; Gao et al., 2021), which requires only a small set of labeled samples and thus reduces annotation costs. Prompt-based learning reformulates the downstream classification task to look more like the masked-language-model task the PLM was trained on. For example, for our task, each input sequence could be appended with a textual prompt, e.g., the propagandistic sequence “Ukraine is full of nazis.” is continued with the prompt “I stand with [MASK]” (which gives the so-called template). Given a mapping of predefined label words to each class (via the so-called verbalizer), the masked language model predicts the probabilities of each label word to fill the [MASK] token and thereby the probabilities of each class. For our task, examples for label words could be “Russia” for the class of propaganda and “Ukraine” for the class of no propaganda. As a result, this introduces the task of prompt engineering, i.e., finding the most suitable template and verbalizer to solve the downstream task. In general, manual prompt engineering can be challenging, especially because the performance in the downstream task depends highly on the prompt (Gao et al., 2021). In our work, we use a three-step procedure: (i) finding the best template, (ii) finding the best verbalizer, and (iii) prompt-based fine-tuning, as follows:

- (i) **Automatic template generation:** Here, we use the LM-BFF procedure from Gao et al. (2021). We randomly sample  $k'$  positive and  $k'$  negative examples for training and validation, which thus requires  $k = 4 \times k'$  samples overall. We use the seq2seq PLM T5 (Raffel et al., 2020) to generate template candidates. Given the training example and an initial verbalizer<sup>5</sup>, T5 then generates a candidate template by filling the missing spans. We use beam search to generate a set of 100 candidate templates. Afterward, we fine-tune each template using the training examples and the downstream PLM. Finally, the best-performing template is chosen based on the performance on the val set.

<sup>4</sup>We also experimented with a variant of the dataset where we used weak labeling for all  $\sim 3.2$  million samples but found comparable results.

<sup>5</sup>This initial verbalizer is only used to generate template candidates. We discard the initial verbalizer for the automatic verbalizer generation in step (ii).

<sup>3</sup>We also evaluated the performance of fine-tuning PLMs with incorporated author and pinned-tweet features. Implementation details and results are reported in Appendix F.

Tweet	BL	PSL
"STOP RUSSIAN AGGRESSION AGAINST #UKRAINE . @USER CLOSE THE SKY OVER UKRAINE ! EXCLUDE RUSSIA FROM THE @USER SECURITY COUNCIL ! #StopPutin #StopRussia HTTPURL"	False	—
"The Textile Worker microdistrict in Donetsk came under fire ! The Ukraine nazis dealt another blow to the residential quarter At least four civilians were killed on the spot . #UkraineRussiaWar #UkraineNazis #ZelenskyWarCriminal @USER @USER HTTPURL"	True	Against Ukraine
"The denazification of Ukraine continues . In Kherson , employees of the Russian Guard detained two accomplices of the Nazis . During the operation , the National Guard officers detained several leaders of neo-Nazi formations and accomplices of the SBU . HTTPURL HTTPURL"	True	Pro Russian government
"Western 'leaders' continue with their irrational drive toward WWII . NATO is a criminal enterprise , an instrument of white power threat to global humanity . Join anti-NATO protests around the world . HTTPURL"	True	Against Western countries
"Chinese and Indian citizens must leave Ukraine because Ukraine is run by the Nazi / Zionist fascists since the coup d'etat of 2014 ."	True	Aimed at other countries

Table 3: Example tweets in our HQP dataset. *BL* is a binary label whether a tweet is propaganda or not. *PSL* is the propaganda strategy.

(ii) **Automatic verbalizer generation:** We use the method from Gao et al. (2021) to generate the verbalizer (i.e., to map predictions to our label classes). For each class, we construct a set of 100 candidate tokens based on the conditional likelihood of the downstream PLM to fill the [MASK] token using the best-performing template from step (i). These candidates are fine-tuned and re-ranked to find the best candidate for each class with regard to the performance on the val set.

(iii) **Prompt-based fine-tuning:** We use the best template from step (i) and the best verbalizer from step (ii) to form our prompt. We fine-tune the downstream PLM with this prompt to create the final model for propaganda detection. We refer to the above model as **LM-BFF**.

In our LM-BFF implementation, we use the OpenPrompt framework (Ding et al., 2022). For template generation, we choose an initial verbalizer with label words “propaganda” (propaganda) and “truth” (no propaganda) and a cloze prompt format (Liu et al., 2023). We choose RoBERTa-large (Liu et al., 2019) as the underlying PLM due to its overall superior performance. We freeze the first 16 layers to control for overfitting and choose a learning rate of  $4e-5$ . We train for a number of 50 epochs and choose the best checkpoint. We set the batch size depending on  $k'$ ; see Table 4. For all other hyper-parameters, we choose the same as those presented in Gao et al. (2021).<sup>6</sup>

<sup>6</sup>We also experimented with prompt-based learning with (i) automatic template generation and a manual verbalizer, (ii) a manual template with automatic verbalizer generation, and (iii) a manual template with a manual verbalizer. However, the results were not better than those reported in Section 5.3. Our explanation for this is the complexity of manually setting up a prompt in our task for propaganda detection.

	$k' = 16$	$k' = 32$	$k' = 64$	$k' = 128$
Batch size	4	8	16	32

Table 4: Batch size of prompt-based learning for different numbers of overall samples ( $k = 4 \times k'$ ).

### 4.3 Extension of LM-BFF to Auxiliary-Task Prompting

We extend the above LM-BFF procedure for inductive learning and use both the *BL* and the *PSL* labels during prompting. The rationale behind this is three-fold: (1) We use information about the propaganda strategy and thus richer labels, which may improve performance. (2) Propaganda can be highly diverse, and, through the use of more granular labels, we can better capture heterogeneity. (3) The overall sample size remains low with only a minor increase in labeling costs. This is beneficial in practice when newly emerging propaganda narratives must be detected and there are thus only a few available samples.

To leverage both *BL* and *PSL* labels, we develop a custom architecture for auxiliary-task prompting, which we refer to as **LM-BFF-AT**. Specifically, we apply steps (i) to (iii) for our two labels *BL* and *PSL*, separately. This results in two different fine-tuned versions of the downstream PLM with different templates and verbalizers. To classify a given input text, we fuse verbalizer probabilities for each label into a classification head, which computes the final prediction. For the classification heads, we train an elastic net and a feed-forward neural network with one hidden layer on top of the verbalizer probabilities. The val set is used for hyper-parameter tuning. Hyper-parameter grids for the classification heads are reported in Appendix E. Note that our LM-BFF-AT approach uses two labels but can

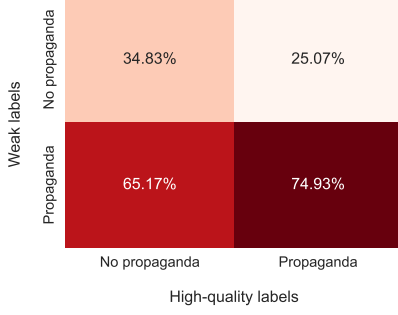


Figure 2: Contingency table comparing weak vs. high-quality labels.

easily be generalized to  $n$  labels.<sup>7</sup>

## 5 Experiments

### 5.1 Discrepancy between Weak Labels vs. High-Quality Labels

**RQ1:** *What is the discrepancy between weak labeling vs. human annotation?*

We find a substantial discrepancy between the weak labels (from HQP-weak) and our high-quality labels (from human annotation). In fact, the overall agreement is only 41.0%. Hence, the majority of labels are thus different. Figure 2 shows the contingency table comparing weak and human annotation. Crucially, this supports our hypotheses that weak labels are noisy and thus often incorrect. This motivates our use of high-quality labels from human annotation in the following.

### 5.2 Propaganda Detection when using Weak vs. High-Quality Labels

**RQ2:** *How well can state-of-the-art PLMs detect online propaganda when trained with weak labels vs. high-quality labels?*

Table 5 compares the performance. For this, we vary the choice of the underlying PLM (BERT, RoBERTa, BERTweet) and what data is used for training (TWE, TWEETSPIN, HQP). We make the following observations: (1) The different PLMs reach a similar performance, which corroborates the robustness and reliability of our results. Recall that we intentionally chose state-of-the-art PLMs to allow for comparability when benchmarking the role of weak vs. high-quality labels. (2) Weak labels from the TWE dataset (Wang et al., 2020) lead to an AUC similar to a random guess, while

<sup>7</sup>We also evaluated the performance of LM-BFF-AT when additionally incorporating author and pinned-tweet features. Implementation details and results are in Appendix F.

weak labels from the TWEETSPIN dataset reach an AUC of 64.03. (3) We use HQP-weak for an ablation study where we use the weak labels from our classification into  $D_+$  and  $D_-$  for training. We register an AUC of 56.79. (4) PLMs trained with high-quality labels perform best with an AUC of 92.25 (for BERTweet). Thereby, we achieve an improvement in AUC over best-performing weak labels (TWEETSPIN) of  $\sim 44\%$ . In sum, the performance gain must be exclusively attributed to the informativeness of high-quality labels (and not other characteristics of the dataset).<sup>8</sup>

We further inspect weak vs. high-quality labels visually. For this, we plot the representation of the [CLS] tokens from HQP using  $t$ -SNE (van der Maaten and Hinton, 2008). As seen in Figure 3, the representations learned with high-quality labels (right plot) are more discriminatory for the true labels than that learned on weak labels (here: TWEETSPIN; left plot).

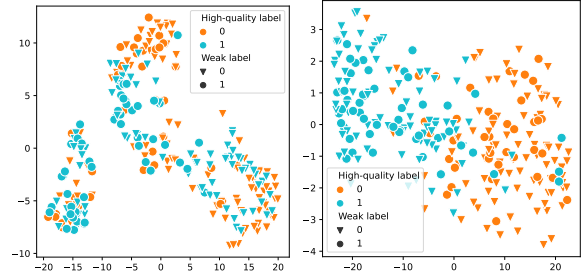


Figure 3: A  $t$ -SNE visualization showing the representations of the [CLS] tokens for BERTweet fine-tuned on TWEETSPIN labels (left) and high-quality labels (right). Evidently, representations learned with high-quality labels are more discriminatory than with weak labels for predictions on the test set.

### 5.3 Performance of Few-Shot Learning

**RQ3:** *How much can few-shot learning reduce labeling costs for detecting online propaganda?*

To reduce the costs of labeling, we further use few-shot learning (see Figure 4). Here, we vary the overall number of labeled samples ( $k = 4 \times k'$ ). We compare the performance of prompt-based learning with LM-BFF (using only  $BL$ ) vs. LM-BFF-AT (using  $BL$  and  $PSL$ ).

<sup>8</sup>Note that the recall improvement with our high-quality labels is relatively small, while we register a strong improvement in precision. In fact, for weak labels, the fine-tuned models tend to predict the propaganda class too often, which leads to a large number of false positives. In practice, this incurs substantial downstream costs during fact-checking (Naumzik and Feuerriegel, 2022) or may oppose free speech rights.

Training data	Precision			Recall			F1			AUC		
	BERT	RoBERTa	BERTweet	BERT	RoBERTa	BERTweet	BERT	RoBERTa	BERTweet	BERT	RoBERTa	BERTweet
TWE (Wang et al., 2020)	14.86 (0.65)	14.67 (0.51)	14.75 (0.20)	46.04 (3.01)	45.71 (2.47)	53.47 (6.13)	22.46 (0.99)	22.20 (0.83)	23.08 (0.52)	47.93 (1.58)	48.63 (1.67)	47.22 (0.81)
TWEETSPIN (Vijayaraghavan and Vosoughi, 2022)	23.08 (1.51)	23.18 (1.11)	23.33 (1.25)	60.09 (1.87)	59.65 (1.48)	59.25 (1.85)	33.32 (1.54)	33.38 (1.36)	33.46 (1.46)	64.03 (1.05)	63.50 (1.06)	63.85 (1.41)
HQP-weak (weak labels on our HQP)	16.42 (0.17)	16.39 (0.31)	16.16 (0.18)	<b>69.24</b> (1.38)	69.94 (3.03)	68.22 (2.64)	26.55 (0.30)	26.56 (0.61)	26.13 (0.43)	56.71 (0.99)	56.79 (2.07)	56.64 (0.76)
<b>HQP (ours)</b>	<b>61.52</b> (5.77)	<b>66.68</b> (2.30)	<b>68.86</b> (2.37)	64.65 (3.83)	<b>70.80</b> (2.85)	<b>70.65</b> (2.52)	<b>62.77</b> (1.92)	<b>68.64</b> (1.80)	<b>69.70</b> (1.31)	<b>88.21</b> (0.62)	<b>91.76</b> (0.62)	<b>92.25</b> (0.80)

Stated: mean (SD).

Table 5: Results of propaganda detection for different PLMs trained on weak vs. high-quality labels.

Figure 4 compares the performance across the different values for  $k'$  and the prompt-based learning methods. Generally, a larger  $k'$  tends to improve the performance. For example, for  $k' = 128$  and LM-BFF, we register a mean F1-score of 43.03 and a mean AUC of 79.74. As expected, this is lower than for fine-tuned PLMs but it is a promising finding since only 2.13% of the labeled examples are used for training and validation. Using LM-BFF-AT with an elastic net as the classification head consistently improves the performance of prompt-based learning across all  $k'$ . For  $k' = 128$ , we achieve a 2.8% improvement in the F1-score (44.22) and a 0.7% improvement in AUC (80.27). On average, over all  $k'$ , the improvement amounts to 1.25% for the F1-score and 0.51% for the AUC. Generally, the variant with an elastic net tends to be better than the variant with a neural network, likely due to the small size of the training sample.<sup>9</sup>

## 6 Discussion

We introduce HQP: the first dataset for online propaganda detection with human annotations. Our experiments further have direct implications for the NLP community.

• **Implication 1:** *When identifying propaganda, there is a substantial discrepancy between weak labeling and human annotations.* This pinpoints weaknesses in existing datasets for online propaganda detection (Wang et al., 2020; Vijayaraghavan and Vosoughi, 2022), since these make exclusive use of weak labeling. To this end, our work highlights the importance of human feedback for sensitive NLP tasks such as propaganda detection.

• **Implication 2:** *High-quality labels are crucial to detect online propaganda.* Our experiments are intentionally based on state-of-the-art PLMs to ensure reliability and comparability of our results.

<sup>9</sup>We report evaluations of the auxiliary task (i.e., prompt-based learning for *PSL*) in Appendix D.

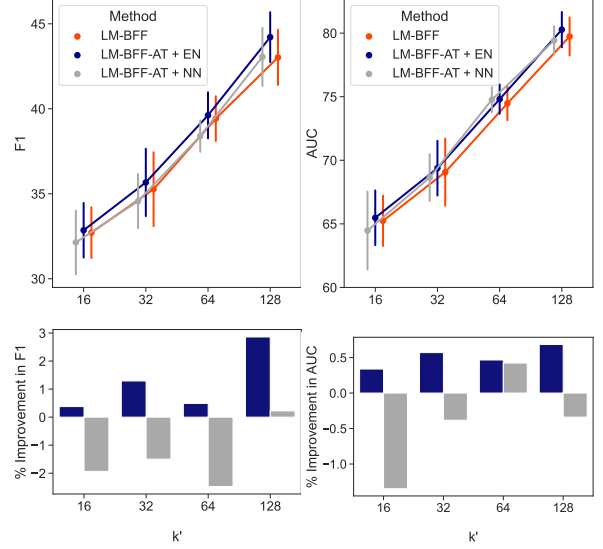


Figure 4: Results for prompt-based learning for LM-BFF vs. LM-BFF-AT (left: F1, right: AUC, top: absolute performance, bottom: %-improvement over LM-BFF). EN (NN) refers to the elastic net (neural net) classification head.  $k'$  refers to the number of examples sampled from each class for both training and validation, while  $k = 4 \times k'$  is the overall number of samples. Error bars denote the standard errors across 5 runs.

Generally, PLMs fail to detect propaganda when fine-tuned with weak labels. In contrast, there is a large improvement ( $\sim 44\%$ ) when using high-quality labels.

• **Implication 3:** *Few-shot learning can be an effective remedy to reduce the cost of human annotation of propaganda.* To the best of our knowledge, our work is the first to adapt few-shot learning (via prompt-based learning) to propaganda detection. Interestingly, our performance is similar to that in related NLP tasks such as, e.g., detecting rumors (Lin et al., 2023) and humor (Li et al., 2023). Despite the challenging nature of our task, the performance of few-shot learning is promising. For example, only  $k = 64$  ( $k' = 16$ ) high-quality annotated samples are needed to outperform propa-



ganda detection with weak labels. For  $k = 512$  ( $k' = 128$ ), we already achieve an improvement over weak labels of 24.54%.

## 7 Limitations

As with other works, ours is not free of limitation. First, there is no universal rule to determine what propaganda is and what not. Hence, the perception may vary across individuals. We address this by having our dataset annotated through multiple raters and showing raters a task description that includes a widely accepted definition of propaganda (see [Smith, 2022](#)). Second, we are further aware that PLMs may embed biases that are populated in downstream tasks. Hence, we call for careful use when deploying our methods in practice. Third, narratives that fall under the scope of propaganda may change over time. Hence, we recommend that both the dataset construction and the PLM fine-tuning is repeated regularly. To this end, we provide a cost-effective approach through few-shot learning.

## 8 Ethics Statement

Our dataset will benefit research on improving social media integrity. The construction of our dataset follows best-practice for ethical research ([Rivers and Lewis, 2014](#)). The dataset construction and usage was approved as ethically unproblematic by the ethics commission of the Faculty of Mathematics, Informatics and Statistics at LMU Munich (ethics approval number: EK-MIS-2023-160). In particular, our dataset contains only publicly available information. The privacy policy of Twitter warns users that their content can be viewed by the general public. Further, we respect the privacy of users and only report aggregate results throughout our paper. Although we believe the intended use of this work is largely positive, there exists potential for misuse (e.g., by propaganda campaigns to run adversarial attacks and develop techniques to avoid detection). To this end, we call for meaningful research by the NLP community to further improve social media integrity. Finally, we encourage careful use of our dataset, as it contains potentially offensive language and manipulative content, which lies in the nature of the task.

## References

Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. 2017. [Deep learning for hate speech detection in tweets](#). In *WWW Companion*.

Alberto Barrón-Cedeño, Israa Jaradat, Giovanni Da San Martino, and Preslav Nakov. 2019. [Proppy: Organizing the news based on their propagandistic content](#). *Information Processing & Management*, 56(5):1849–1864.

Tian Bian, Xi Xiao, Tingyang Xu, Peilin Zhao, Wenbing Huang, Yu Rong, and Junzhou Huang. 2020. [Rumor detection on social media with bi-directional graph convolutional networks](#). In *AAAI*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). In *NeurIPS*.

Jacob Cohen. 1960. [A coefficient of agreement for nominal scales](#). *Educational and Psychological Measurement*, 20(1):37–46.

Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, Rostislav Petrov, and Preslav Nakov. 2019. [Fine-grained analysis of propaganda in news articles](#). In *EMNLP-IJCNLP*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *NAACL*.

Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Haitao Zheng, and Maosong Sun. 2022. [OpenPrompt: An open-source framework for prompt-learning](#). In *ACL*.

Antigoni Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. [Large scale crowdsourcing and characterization of Twitter abusive behavior](#). In *ICWSM*.

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *ACL-IJCNLP*.

Dominique Geissler, Dominik Bär, Nicolas Pröllochs, and Stefan Feuerriegel. 2022. [Russian propaganda on social media during the 2022 invasion of Ukraine](#). arXiv.

Yiqiao Jin, Xiting Wang, Ruichao Yang, Yizhou Sun, Wei Wang, Hao Liao, and Xing Xie. 2022. [Towards fine-grained reasoning for fake news detection](#). In *AAAI*.

Adam Kowalski. 2022. [Disinformation and Russia’s war of aggression against Ukraine: Threats and governance responses](#). *OECD Report*.

Klaus Krippendorff, Yann Mathet, Stéphane Bouvry, and Antoine Widlöcher. 2016. [On the reliability of unitizing textual continua: Further developments](#). *Quality & Quantity*, 50(6):2347–2364.

- Valerio La Gatta, Chiyu Wei, Luca Luceri, Francesco Pierri, and Emilio Ferrara. 2023. [Retrieving false claims on Twitter during the Russia-Ukraine conflict](#). In *WWW Companion*.
- Junze Li, Mengjie Zhao, Yubo Xie, Antonis Maronikolakis, Pearl Pu, and Hinrich Schütze. 2023. [This joke is \[MASK\]: Recognizing humor and offense with prompting](#). *Transfer Learning for Natural Language Processing Workshop*.
- Hongzhan Lin, Pengyao Yi, Jing Ma, Haiyun Jiang, Ziyang Luo, Shuming Shi, and Ruifang Liu. 2023. [Zero-shot rumor detection with propagation structure via prompt learning](#). In *AAAI*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *ACM Computing Surveys*, 55(9):1–35.
- Yinhan Liu, Myle Ott, Naman Goyal, Du Jingfei, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A robustly optimized BERT pretraining approach](#). arXiv.
- Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020. [Fine-grained fact verification with kernel graph attention network](#). In *ACL*.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *ICLR*.
- Yi-Ju Lu and Cheng-Te Li. 2020. [GCAN: Graph-aware co-attention networks for explainable fake news detection on social media](#). In *ACL*.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. [HateXplain: A benchmark dataset for explainable hate speech detection](#). In *AAAI*.
- Christof Naumzik and Stefan Feuerriegel. 2022. [Detecting false rumors from retweet dynamics on social media](#). In *WWW*.
- Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen. 2020. [BERTweet: A pre-trained language model for English tweets](#). In *EMNLP*.
- John Pavlopoulos, Leo Laugier, Alexandros Xenos, Jeffrey Sorensen, and Ion Androutsopoulos. 2022. [From the detection of toxic spans in online discussions to the analysis of toxic-to-civil transfer](#). In *ACL*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, and Dario Amodei. 2019. [Language models are unsupervised multitask learners](#). *Preprint*.
- Collin 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](#). *JMLR*, 21(140):1–67.
- Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. [Truth of varying shades: Analyzing language in fake news and political fact-checking](#). In *EMNLP*.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using siamese BERT-networks](#). In *EMNLP-IJCNLP*.
- Caitlin M. Rivers and Bryan L. Lewis. 2014. [Ethical research standards in a world of big data](#). *F1000Research*, 3(8).
- Timo Schick and Hinrich Schütze. 2021. [It’s not just size that matters: Small language models are also few-shot learners](#). In *NAACL*.
- Bruce Lannes Smith. 2022. [Propaganda](#). *Encyclopedia Britannica*.
- Veronika Solopova, Oana-Iuliana Popescu, Christoph Benz Müller, and Tim Landgraf. 2023. [Automated multilingual detection of Pro-Kremlin propaganda in newspapers and Telegram posts](#). arXiv.
- Hyunjin Song, Petro Tolochko, Jakob-Moritz Eberl, Olga Eisele, Esther Greussing, Tobias Heidenreich, Fabienne Lind, Sebastian Galyga, and Hajo G. Boomgaarden. 2020. [In validations we trust? The impact of imperfect human annotations as a gold standard on the quality of validation of automated content analysis](#). *Political Communication*, 37(4):550–572.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: A large-scale dataset for Fact Extraction and VERification](#). In *NAACL*.
- Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-SNE](#). *JMLR*, 9:2579–2605.
- Prashanth Vijayaraghavan and Soroush Vosoughi. 2022. [TWEETSPIN: Fine-grained propaganda detection in social media using multi-view representations](#). In *NAACL*.
- Liqiang Wang, Xiaoyu Shen, Gerard de Melo, and Gerhard Weikum. 2020. [Cross-domain learning for classifying propaganda in online contents](#). In *TTO*.
- Lingwei Wei, Dou Hu, Wei Zhou, Zhaojuan Yue, and Songlin Hu. 2021. [Towards propagation uncertainty: Edge-enhanced bayesian graph convolutional networks for rumor detection](#). In *ACL-IJCNLP*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *EMNLP*.

Rui Xia, Kaizhou Xuan, and Jianfei Yu. 2020. [A state-independent and time-evolving network for early rumor detection in social media](#). In *EMNLP*.

Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. [Defending against neural fake news](#). In *NeurIPS*.

Kaimin Zhou, Chang Shu, Binyang Li, and Jey Han Lau. 2019. [Early rumour detection](#). In *NAACL*.

## A Keywords for Dataset Construction

Table A.1 lists the keywords used in our dataset construction process to obtain candidate tweets for the positive class ( $D_+$ ) via accusations in replies. Table A.2 shows the keywords that are used to collect candidate tweets for the negative class ( $D_-$ ). Generally, keywords relevant to the positive class should mostly be terms that express accusations of propaganda, while keywords relevant to the negative class should be mostly terms that refer to general activities of the war. The keywords in both lists contain further results from our construction procedure in that we list the iteration in which they were added to the list.

Keywords ( $D_+$ )	Iteration
russia(n) $\wedge$ propaganda	1
russia(n) $\wedge$ propagandist	1
kremlin $\wedge$ propaganda	2
kremlin $\wedge$ propagandist	2
putinist(s)	2
putinism	2
russia(n) $\wedge$ lie(s)	3
war $\wedge$ propaganda	3
war $\wedge$ lie(s)	3
putin $\wedge$ propaganda	3
putin $\wedge$ propagandist	3
russia(n) $\wedge$ fake news	3

Table A.1: List of keywords used to get propaganda accusations and the corresponding iteration they were added to the keyword list. The  $\wedge$ -operator indicates that both keywords have to appear.

Keywords ( $D_-$ )	Iteration
russia $\wedge$ war	1
ukraine $\wedge$ war	1
#istandwithrussia	1
#istandwithputin	1
russian $\wedge$ war	2
ukrainian $\wedge$ war	2
#russianukrainianwar	2
#ukrainerussiawar	2
#standwithrussia	2
#standwithputin	2
#russia	2
#russiaukraine	2
#ukraine	2
putin $\wedge$ war	3
#putin	3
#lavrov	3
#zakharova	3
#nato	3
#donbass	3
#mariupol	3

Table A.2: List of keywords used to get war-related tweets and the corresponding iteration they were added to the keyword list. The  $\wedge$ -operator indicates that both keywords have to appear.

## B Annotation Instructions

Figure B.1 shows the instructions of batch annotations we present to the workers on Prolific (<https://www.prolific.co/>). We follow best practices (Song et al., 2020). That is, we provide a detailed and comprehensible description of the task, a precise definition of the labels, and a transparent disclosure that we use attention checks.

Given the complexity and subjective nature of propaganda, we put great emphasis on providing meaningful and precise instructions for annotation. Therefore, the final instructions in Figure B.1 are a result of several iterations of improvement, each followed by an internal discussion and analysis. In each iteration, we had three workers from Prolific annotate a random sample of 100 tweets using the current version of the instructions. We analyzed the resulting annotations and focused on the samples with disagreement. We then aimed to address the corresponding issues in the next update of our instructions. In each iteration, we calculated the inter-annotator agreement using Krippendorff’s alpha (Krippendorff et al., 2016) and stopped iterating after surpassing an agreement of 0.8 for the first time. As a result, we updated the instructions three times.

We also followed best practices and applied attention checks. Participants failing the attention checks were removed, and the annotation was repeated.



## Instructions

First, we would like to thank you for your interest on helping with Tweet annotations regarding their propagandistic content. Please read the following instructions carefully.

**Definition:** Please use the following definition of Propaganda to guide you through the annotation process:

*"Propaganda is expression of opinion or action by individuals or groups deliberately designed to influence opinions or actions of other individuals or groups with reference to predetermined ends."*

**Propaganda Detection:** Your task will be to annotate 100 Tweets regarding their propagandistic content. All Tweets included in this task relate to the Russian-Ukrainian conflict. If you feel like you need more background information on the Russian-Ukrainian conflict, we provide you with the following articles:

- Russian government accounts are using a Twitter loophole to spread disinformation (theconversation.com)
- Russia is swaying Twitter users outside the West to its side - The Economist (economist.com)
- Russia Takes Censorship to New Extremes, Stifling War Coverage - The New York Times (nytimes.com)
- Key Moments in the Russia-Ukraine War: A Timeline - The New York Times (nytimes.com)

In your collection of Tweets, some of them will refer to certain conflict related events or facts. To know whether these Tweets contain Russian misinformation (which is a common propagandistic strategy), we ask you to perform a quick fact checking, if needed.

To further help you with the annotations, we provide the following list of notes on what to regard as propaganda and what not:

- Misinformation in favour of the Russian government is regarded as propaganda
- Propaganda which is not in favour of the Russian government is not regarded as propaganda here (we are only interested in detecting Russian propaganda)
- If only small parts of the Tweet contains propagandistic content, we regard the entire tweet as propagandistic
- Tweets containing a website/image/video URL and where the Tweet content itself does not classify into propaganda or not, we regard as non-propagandistic
- Tweets designed to spread pro-Russian-government stance (in form of e.g. slogans, hashtags, ...) are also regarded as propaganda

**Propaganda Strategy:** In addition, we ask you to annotate the strategy behind the propagandistic content. Russian propaganda is known to influence opinions around the world not only regarding the Russian government. If you decided that a specific tweet contains propaganda, please specify whether it is designed to influence opinions regarding:

- against Western countries
- against Ukraine
- pro Russian government
- aimed at other countries

**Attention Checks:** Please note that we have included attention checks to ensure reliable annotations. On the following page we will ask you some basic questions regarding these instructions. If you have read these instructions carefully, you will be able to answer them all correctly. Furthermore please note, that during the attention check you can go back to the instructions and find the answers to the questions. If your answers to these questions are incorrect, the survey will redirect you to Prolific and you will not be rewarded. Furthermore, we included synthetic Tweets, where the classification into propaganda or no-propaganda is obvious.

Figure B.1: Instructions for annotators.

## C Implemenation Details for Full Fine-tuning

For fine-tuning, we add a linear layer to the hidden representation of the [CLS] token. The PLMs are then fine-tuned using the transformer framework from Huggingface (Wolf et al., 2020). We set the maximum sequence length to 128. We use a training batch size of 32 and a learning rate of  $4e-5$ . We freeze the first 16 layers of the PLMs. For BERT-large and RoBERTa-large, we add emoji-tokens to the vocabulary due to their frequent and meaningful use in social media.<sup>10</sup> Weight updates are performed using the AdamW-optimizer (Loshchilov and Hutter, 2019). We fine-tune for a maximum number of 5 epochs. We validate the performance every 500 steps for fine-tuning with TWEETSPIN and every 50 steps otherwise. Early stopping is used when the loss on the validation set does not decrease for more than 5 validation steps.

## D Evaluation of the Auxiliary Task in LM-BFF-AT

In Table D.1, we report the performance of the auxiliary task in LM-BFF-AT, i.e., the performance of prompt-based learning using *PSL*.

k'	Weighted Precision	Weighted Recall	Weighted F1
16	75.06 (2.54)	66.17 (6.81)	69.84 (4.02)
32	75.6 (1.21)	64.93 (10.73)	68.75 (7.48)
64	77.46 (2.19)	75.31 (4.45)	76.13 (2.15)
128	77.18 (3.45)	74.39 (3.87)	75.41 (1.86)

Stated: mean (SD).

Table D.1: Evaluation results for the auxiliary task of prompt-based learning using *PSL* on HQP.

## E Hyper-Parameters for Classification Heads of LM-BFF-AT

For our extension of the LM-BFF method, namely LM-BFF-AT, we perform hyper-parameter tuning using grid search for the two classification heads, i.e., the elastic net and the neural net. The tuning grids are reported in Table E.1. We implement the elastic net using Python’s scikit-learn module. The neural net is implemented using PyTorch.

<sup>10</sup>For BERTweet-large, emoji-tokens were already incorporated in the vocabulary during training.

Classification head	Hyper-parameter	Grid
Elastic net	Cost	{0.1, 0.25, 0.5, 1, 2, 4, 8}
	L1-ratio	{0.1, 0.15, 0.2, 0.25, 0.3, 0.35}
Neural net	Dropout	{0.2, 0.4, 0.6}
	Learning rate	{0.001, 0.01, 0.02}
	Batch size	{2, 4, 8, 16}
	Neurons in hidden layer	$input\_dim * \{0.5, 1, 1.5, 2\}$

Table E.1: Grids for hyper-parameter tuning of the two classification heads for prompt-based learning with LM-BFF-AT. The cost  $C = \frac{1}{\lambda}$  is the inverse of the regularization strength. L1-ratio is the elastic net mixing parameter (i.e., is equivalent to only using an L2-penalty (or L1-penalty) for L1-ratio= 0 (or L1-ratio= 1) and mixes both for  $0 < \text{L1-ratio} < 1$ ). Here,  $input\_dim$  is the number of input features and therefore depends on whether only probabilities from the verbalizers or also author representations are used.

## F Propaganda Detection with Additional Meta Information

We extend our propaganda detection so that we not only use the content but the additional meta information (i.e., author features and pinned-tweet features) for propaganda detection.

### F.1 Data Enrichment

We enrich HQP with additional meta information from the social network. Here, we use a comprehensive set of author features (e.g., number of followers, account age, verified status) and pinned-tweet<sup>11</sup> features (e.g., tweet age, number of likes, number of retweets). Our data further includes the profile description of authors using embeddings from SBERT (Reimers and Gurevych, 2019). We choose SBERT due to its strength in capturing meaningful representations from short text (Reimers and Gurevych, 2019). A few tweets corresponded to authors whose accounts were already deleted, which reduces our dataset to the final size of  $N = 29,596$ . We note that for fair benchmarking we used the dataset of size  $N = 29,596$  for all experiments. The full list of features is in Appendix F.2.

### F.2 Summary Statistics

Table F.1 reports summary statistics of the author features for HQP. We compare mean and standard

<sup>11</sup>Every Twitter user can choose to pin one (self-written) tweet to her/his account, which is then always displayed at the top of the profile.

	Propaganda = true		Propaganda = false	
	Mean	Median	Mean	Median
Verified (=1; not=0)	0.09	0	0.19	0
#Followers	102784.72	513	595308.01	13405
#Following	1927.58	530.5	3237.82	786
#Tweets	37780.17	11123.5	64060.48	16811
#Listed*	633.36	5	2449.10	11
Account age (in days)	2272.17	1938	2763.28	2895
#Followers divided by Account age	28.05	0.37	161.74	0.71
#Following divided by Account age	1.19	0.34	1.48	0.34
#Tweets divided by Account age	19.34	7.915	25.76	8.58

\*Number of Twitter lists (i.e., a curated group of accounts) comprising the account of the author.

Table F.1: Summary statistics of author features in HQP.

	Propaganda = true		Propaganda = false	
	Mean	Median	Mean	Median
Account age (in days)	365.30	198	366.82	170
#Retweets	573.12	13	1051.95	17
#Replies	251.85	6	416.71	7
#Likes	2232.50	39	4336.85	62
#Quotes	70.19	2	131.85	2
#Retweets divided by account age	16.70	0.08	30.49	0.14
#Replies divided by account age	5.09	0.03	11.95	0.05
#Likes divided by account age	59.17	0.25	111.47	0.5
#Quotes divided by account age	1.14	0.01	3.18	0.01

Table F.2: Summary statistics for pinned-tweet features of authors in HQP.

deviation separately for both propagandistic and non-propagandistic content. In line with previous findings (Geissler et al., 2022), we find that authors of propagandistic content are, on average, less often verified, have fewer followers and tweets, and are characterized by a younger account age. Table F.2 reports summary statistics for the pinned-tweet features of the authors in HQP.

### F.3 Adaptation of Methods

**(i) Full fine-tuning:** We modify the classification head so that we perform full fine-tuning with author and pinned-tweet features. The self-description of tweet authors is encoded into a 768-dimensional vector using SBERT. We normalize the numerical features from Table F.1 and Table F.2 and append them to the self-description vector to get an author representation. We concatenate the hidden representation of the [CLS] token generated by the PLM and the author representation and again feed them to a linear layer. The fine-tuning procedure and the hyper-parameters are identical to those described in Section 4.1.

**(ii) Prompt-based learning:** We create the same author representation as in (i). Here, we concatenate the author representation with the verbalizer probabilities for both *BL* and *PSL* and again feed them to a classification head as in LM-BFF-AT. The prompt-based learning procedure is identical to

that in Section 4.2, i.e., we select the optimal template and verbalizer from each run in Section 5.3. The tuning grids for hyper-parameter tuning for the two classification heads are identical to those in Table E.1.

### F.4 Results

Table F.3 reports the prediction performance when additionally using author features. Values are in bold if the model that uses additional author and pinned-tweet features outperforms the counterpart without author features. Overall, we observe a tendency that the results improve when additionally using author features. This tendency is seen for both full fine-tuning and prompt-based learning. However, the performance gain from using the content is larger than the performance gain from using author features. This can be expected as propaganda spreaders typically do not explicitly disclose their manipulative intention but instead aim to deceive users.

	PLM	P	R	F1	AUC
(i) Full fine-tuning	BERT	<b>64.58</b> (3.19)	62.40 (1.94)	<b>63.42</b> (1.64)	<b>88.54</b> (0.81)
	RoBERTa	<b>67.61</b> (3.05)	69.59 (3.02)	68.52 (2.04)	91.37 (1.02)
	BERTweet	67.06 (3.39)	<b>72.68</b> (1.54)	<b>69.71</b> (1.79)	92.06 (0.86)
	$k'$	P	R	F1	AUC
(ii) Prompt-based learning	16	22.38 (3.37)	<b>64.94</b> (9.79)	<b>32.98</b> (3.49)	64.99 (4.16)
	32	<b>24.75</b> (6.24)	65.16 (3.6)	<b>35.70</b> (4.16)	69.29 (4.48)
	64	<b>27.73</b> (2.80)	70.72 (4.07)	<b>39.74</b> (2.88)	<b>74.88</b> (2.21)
	128	30.81 (2.86)	76.38 (1.52)	43.86 (3.02)	79.46 (2.89)

Stated: mean (SD). P: precision. R: recall.

Table F.3: Evaluation results of (i) full fine-tuning and (ii) prompt-based learning on HQP while incorporating author and pinned-tweet features. Results are shown in **bold** if the performance with author and pinned-tweet features is better than the performance without author features.