Overview of the WANLP 2022 Shared Task on Propaganda Detection in Arabic

Firoj Alam¹, Hamdy Mubarak¹, Wajdi Zaghouani², Giovanni Da San Martino³, Preslav Nakov⁴

¹Qatar Computing Research Institute, HBKU, Qatar

²Hamad Bin Khalifa University, Qatar

³University of Padova, Italy

⁴Mohamed bin Zayed University of Artificial Intelligence, UAE

{falam,hmubarak,wzaghouani}@hbku.edu.qa,dasan@math.unipd.it,preslav.nakov@mbzuai.ac.ae

Abstract

Propaganda is the expression of an opinion or an action by an individual or a group deliberately designed to influence the opinions or the actions of other individuals or groups with reference to predetermined ends, which is achieved by means of well-defined rhetorical and psychological devices. Propaganda techniques are commonly used in social media to manipulate or to mislead users. Thus, there has been a lot of recent research on automatic detection of propaganda techniques in text as well as in memes. However, so far the focus has been primarily on English. With the aim to bridge this language gap, we ran a shared task on detecting propaganda techniques in Arabic tweets as part of the WANLP 2022 workshop, which included two subtasks. Subtask 1 asks to identify the set of propaganda techniques used in a tweet, which is a multilabel classification problem, while Subtask 2 asks to detect the propaganda techniques used in a tweet together with the exact span(s) of text in which each propaganda technique appears. The task attracted 63 team registrations, and eventually 14 and 3 teams made submissions for subtask 1 and 2, respectively. Finally, 11 teams submitted system description papers.

1 Introduction

Social media platforms have become an important communication channel, where we can share and access information from a variety of sources. Unfortunately, the rise of this democratic information ecosystem was accompanied by and dangerously polluted with misinformation, disinformation, and malinformation in the form of propaganda, conspiracies, rumors, hoaxes, fake news, hyper-partisan content, falsehoods, hate speech, cyberbullying, etc. (Oshikawa et al., 2020; Alam et al., 2021; Pramanick et al., 2021; Rosenthal et al., 2021; Alam et al., 2022; Barnabò et al., 2022; Guo et al., 2022; Hardalov et al., 2022; Nguyen et al., 2022; Sharma et al., 2022)

Propaganda is conveyed through the use of diverse propaganda techniques (Miller, 1939), which range from leveraging on the emotions of the audience (e.g., using loaded language, appealing to fear, etc.) to using logical fallacies such as straw men (misrepresenting someone's opinion), whataboutism, red herring (presenting irrelevant data), etc. In the last decades, propaganda was widely used on social media to influence and/or mislead the audience, which became a major concern for different stakeholders, social media platforms, and policymakers. To address this problem, the research area of computational propaganda has emerged, and here we are particularly interested in automatically identifying the use of propaganda techniques in text, images, and multimodal content. Prior work in this direction includes identifying propagandistic content in an article based on writing style and readability level (Rashkin et al., 2017; Barrón-Cedeno et al., 2019), at the sentence and the fragment levels from news articles with finegrained techniques (Da San Martino et al., 2019b), and in memes (Dimitrov et al., 2021a). These efforts focused on English, and there was no prior work on Arabic. Our shared task aims to bridge this gap by focusing on detecting propaganda in Arabic social media text, i.e., tweets.

2 Related Work

In the current information ecosystem, propaganda has evolved to *computational propaganda* (Woolley and Howard, 2018; Da San Martino et al., 2020b), where information is distributed on social media platforms, which makes it possible for malicious users to reach well-targeted communities at high velocity. Thus, research on propaganda detection has focused on analyzing not only news articles but also social media content (Rashkin et al., 2017; Barrón-Cedeno et al., 2019; Da San Martino et al., 2019b, 2020b; Nakov et al., 2021a,b; Hristakieva et al., 2022).

Rashkin et al. (2017) focused on article-level propaganda analysis. They developed the TSHP-17 corpus, which used distant supervision for annotation with four classes: trusted, satire, hoax, and propaganda. The assumption of their distant supervision approach was that all articles from a given news source should share the same label. They collected their articles from the English Gigaword corpus and from seven other unreliable news sources, including two propagandistic ones. Later, Barrón-Cedeno et al. (2019) developed a new corpus, QProp, with two labels: propaganda vs. non-propaganda, and also experimented on TSHP-17 and QProp corpora. For the TSHP-17 corpus, they binarized the labels: propaganda vs. any of the other three categories as non-propaganda. They investigated the writing style and the readability level of the target document, and trained models using logistic regression and SVMs. Their findings confirmed that using distant supervision, in conjunction with rich representations, might encourage the model to predict the source of the article, rather than to discriminate propaganda from non-propaganda. Similarly, Habernal et al. (2017, 2018) developed a corpus with 1.3k arguments annotated with five fallacies, including ad hominem, red herring, and irrelevant authority, which directly relate to propaganda techniques.

Recently, Da San Martino et al. (2019b), curated a set of persuasive techniques, ranging from leveraging on the emotions of the audience such as using loaded language and appeal to fear, to logical fallacies such as straw man (misrepresenting someone's opinion) and red herring (presenting irrelevant data). They focused on textual content, i.e., newspaper articles. In particular, they developed a corpus of news articles annotated with eighteen propaganda techniques. The annotation was at the fragment level, and could be used for two tasks: (i) binary classification —given a sentence in an article, predict whether any of the 18 techniques has been used in it, and (ii) multi-label classification and span detection task —given a raw text, identify both the specific text fragments where a propaganda technique is used as well as the specific technique. They further proposed a multigranular deep neural network that captures signals from the sentence-level task and helps to improve the fragment-level classifier. Da San Martino et al. (2020a) also organized a shared task on Detection of Propaganda Techniques in News Articles.

Subsequently, Dimitrov et al. (2021b) organized the SemEval-2021 task 6 on Detection of Propaganda Techniques in Memes. It had a multimodal setup, combining text and images, and asked participants to build systems to identify the propaganda techniques used in a given meme. Yu et al. (2021) looked into interpretable propaganda detection.

Other related shared tasks include the FEVER task (Thorne et al., 2018) on fact extraction and verification, the Fake News Challenge (Hanselowski et al., 2018), the FakeNews task at MediaEval (Pogorelov et al., 2020), as well as the NLP4IF tasks on propaganda detection (Da San Martino et al., 2019a) and on fighting the COVID-19 infodemic in social media (Shaar et al., 2021a). Finally, we should mention the CheckThat! lab at CLEF (Elsayed et al., 2019a,b; Barrón-Cedeño et al., 2020; Shaar et al., 2020; Hasanain et al., 2020; Nakov et al., 2021c,d; Shaar et al., 2021b; Nakov et al., 2022a,b,c,d), which addresses many aspects of disinformation for different languages over the years such as fact-checking, verifiable factual claims, check-worthiness, attentionworthiness, and fake news detection.

The present shared task is inspired from prior work on propaganda detection. In particular, we adapted the annotation instructions and the propaganda techniques discussed in (Da San Martino et al., 2019b; Dimitrov et al., 2021b).

3 Tasks and Dataset

Below, we first formulate the two subtasks of our shared task, and then we discuss our datasets, including how we collected the data and what annotation guidelines we used.

3.1 Tasks

In the shared tasks, we offered the following two subtasks:

- **Subtask 1:** Given the text of a tweet, identify the propaganda techniques used in it.
- **Subtask 2:** Given the text of a tweet, identify the propaganda techniques used in it together with the span(s) of text in which each propaganda technique appears.

Note that Subtask 1 is formulated as a multilabel classification problem, while Subtask 2 is a sequence labeling task.

Figure 1: An example of tweet annotation with propaganda techniques loaded language and name calling.

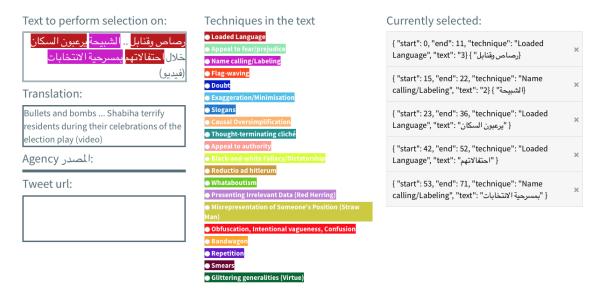
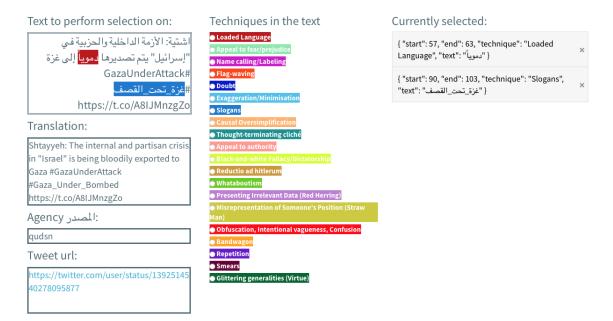


Figure 2: An example of tweet annotation with propaganda techniques loaded language and slogan.



3.2 Dataset

We used Social Bakers¹ to obtain the top-2 news sources from each Arab country, e.g., Al Arabiya and Sky News Arabia from UAE, Al Jazeera and Al Sharq from Qatar, etc. We further added five international sources that broadcast Arabic news: Al-Hurra News, BBC Arabic, CNN Arabic, France 24, and Russia Today. We then extracted from Twitter their latest 3,200 tweets. To have a balanced dataset that covers a wide range of topics, we chose 100 random tweets from each source, and then we sampled 930 tweets for annotation.

We target emotional appeals (e.g., loaded language, appeal to fear, flag waving, exaggeration, etc.) and logical fallacies (e.g., whataboutism, causal oversimplification, red herring, band wagon, etc.). We adopted the same techniques studied in (Da San Martino et al., 2019b; Dimitrov et al., 2021b). Below we briefly summarize them:

1. **Appeal to authority:** Stating that a claim is true simply because a valid authority or expert on the issue said it was true. We also include here the special case where the reference is not an authority or an expert, which is referred to as *Testimonial* in the literature.

¹https://www.socialbakers.com/

- Appeal to fear / prejudices: Seeking to build support for an idea by instilling anxiety and/or panic in the population towards an alternative. In some cases, the support is built based on preconceived judgements.
- Bandwagon Attempting to persuade the target audience to join in and take the course of action because "everyone else is taking the same action."
- 4. **Black-and-white fallacy or dictatorship:** Presenting two alternative options as the only possibilities, when in fact more possibilities exist. As an the extreme case, tell the audience exactly what actions to take, eliminating any other possible choices (ictatorship).
- 5. Causal oversimplification: Assuming a single cause or reason when there are actually multiple causes for an issue. This includes transferring blame to one person or group of people without investigating the complexities of the issue.
- 6. **Doubt:** Questioning the credibility of someone or something.
- 7. **Exaggeration / minimisation:** Either representing something in an excessive manner: making things larger, better, worse (e.g., *the best of the best, quality guaranteed*) or making something seem less important or smaller than it really is (e.g., saying that an insult was actually just a joke).
- 8. **Flag-waving:** Playing on strong national feeling (or to any group, e.g., race, gender, political preference) to justify or to promote an action or an idea.
- 9. Glittering generalities (virtue) These are words or symbols in the value system of the target audience that produce a positive image when attached to a person or issue. Peace, hope, happiness, security, wise leadership, freedom, "The Truth", etc. are virtue words. Virtue can be also expressed in images, where a person or an object is depicted positively.
- Loaded language: Using specific words and phrases with strong emotional implications (either positive or negative) to influence an audience.

- 11. **Misrepresentation of someone's position** (**straw man**): Substituting an opponent's proposition with a similar one, which is then refuted in place of the original proposition.
- 12. **Name calling or labeling:** Labeling the object of the propaganda campaign as something that the target audience fears, hates, finds undesirable or loves, praises.
- 13. **Obfuscation, intentional vagueness, confusion:** Using words that are deliberately not clear, so that the audience may have their own interpretations. For example, when an unclear phrase with multiple possible meanings is used within an argument and, therefore, it does not support the conclusion.
- 14. **Presenting irrelevant data (red herring):** Introducing irrelevant material to the issue being discussed, so that everyone's attention is diverted away from the points made.
- 15. **Reductio ad hitlerum:** Persuading an audience to disapprove an action or an idea by suggesting that the idea is popular with groups hated in contempt by the target audience. It can refer to any person or concept with a negative connotation.
- 16. **Repetition:** Repeating the same message over and over again, so that the audience will eventually accept it.
- 17. **Slogans:** A brief and striking phrase that may include labeling and stereotyping. Slogans tend to act as emotional appeals.
- 18. **Smears** A smear is an effort to damage or call into question someone's reputation, by propounding negative propaganda. It can be applied to individuals or groups.
- 19. **Thought-terminating cliché:** Words or phrases that discourage critical thought and meaningful discussion about a given topic. They are typically short, generic sentences that offer seemingly simple answers to complex questions or that distract the attention away from other lines of thought.
- 20. **Whataboutism:** A technique that attempts to discredit an opponent's position by charging them with hypocrisy without directly disproving their argument.

Table 1: Statistics about the corpus. In parentheses, we show the number of tweets. *Total* represents the number of techniques in each set.

Prop Technique	Train (504)		Dev-Test (51)	Test (323)	
Appeal to authority	21	7	1	1	
Appeal to fear/prejudice	48	7	4	25	
Black-and-white Fallacy/Dictatorship	2	1	2	7	
Causal Oversimplification	4	1	1	4	
Doubt	29	1	2	19	
Exaggeration/Minimisation	44	10	16	26	
Flag-waving	5	2	2	9	
Glittering generalities	25	7	2	1	
(Virtue)	25	/	2	1	
Loaded Language	446	46	42	326	
Name calling/Labeling	244	44	33	163	
Obfuscation, Intentional vagueness, Confusion	9	3	1	6	
Presenting Irrelevant Data (Red Herring)	1	0	0	0	
Repetition	9	2	1	3	
Slogans	44	1	1	6	
Smears	85	12	15	50	
Thought-terminating cliché	6	1	1	0	
Whataboutism	3	1	1	0	
Total	1025	146	125	646	

The annotation is done in different stages: (i) three annotators independently annotate the same tweet, and (ii) they meet together with one consolidator to discuss each instance and to come up with gold annotations. Since the annotations are at the fragment level, it might happen that an annotation is spotted by only one annotator. The two phases ensure that each annotation is eventually discussed by all annotators. In order to train the annotators, we provide clear annotation instructions with examples and ask them to annotate a sample of tweets. Then, we revise their annotations and provide feedback. Figures 1 and 2 show example tweets with annotated propaganda techniques.

Table 1 shows the distribution of the propaganda techniques in our dataset for different data splits. Our annotation guidelines inclide twenty techniques, but in the annotated dataset, there were no instances of *bandwagon*, *straw man*, and *reductio ad hitlerum*. Overall, the distribution of the propaganda techniques in our dataset is very skewed, which made the task challenging.

4 Evaluation Framework

4.1 Evaluation Measures

To measure the performance of the systems, for both subtasks, we use micro-F1 and macro-F1, as these are multi-class multi-label problems, where the labels are imbalanced. The official evaluation measure for subtask 1 is micro-F1, but the scorer also reports macro-F1.

Subtask 2 is a multi-label sequence tagging problem. We modified the standard micro-averaged F1 to account for partial matching between the spans. More details about the modified macro-averaged F1 can be found in (Da San Martino et al., 2019b; Dimitrov et al., 2021b).

4.2 Task Organization

We ran the shared task in two phases:

Development Phase In the first phase, we provided the participants three subsets of the dataset: train, dev, and dev_test. The purpose of the dev set was to fine-tune the trained model, and the dev_test set was to evaluate the model performance on unseen dev_test set.

Test Phase In the second phase, we released the actual test set and the participants were given just a few days to submit their final predictions via the submission system on Codalab.² In this phase, the participants could again submit multiple runs, but they would not get any feedback on their performance. Only the latest submission of each team was considered as official and was used for the final team ranking. The final leaderboard on the test set was made publicly available after the system submission deadline.

5 Participants and Results

In this section, we provide a general description of the systems that participated in each subtask and their results. Table 2 shows the results for all teams for both subtasks, as well as a random baseline. We can see that subtask 1 was more popular, attracting submissions by 14 teams, while there were only three submissions for subtask 2.

5.1 Subtask 1

Table 3 gives an overview of the systems that took part in subtask 1. We can see that transformers were quite popular, most notably AraBERT, followed by BERT, and MARBERT. Some participants also used ensembles methods, data augmentation, and standard preprocessing.

The best-performing team NGU_CNLP (Samir et al., 2022) first explored various baselines models such as bag of words with SVM, Naïve Bayes, Stochastic Gradient Descent, Logistic Regression,

²https://codalab.lisn.upsaclay.fr/ competitions/7274

Table 2: Results for subtask 1 on multilabel propaganda detection and subtask 2 on identifying propaganda techniques and their span(s) in the text. The results are ordered by the official score: Micro-F1. *Indicated that no system description paper was submitted.

Rank/Team	Macro F1	Micro F1					
Subtask 1							
1. NGU_CNLP (Samir et al., 2022)	0.185	0.649					
2. IITD (Mittal and Nakov, 2022)	0.183	0.609					
3. CNLP-NITS-PP (Laskar et al., 2022)	0.068	0.602					
3. AraBEM (Eshrag Ali et al., 2022)	0.068	0.602					
3. Pythoneers (Attieh and Hassan, 2022)	0.177	0.602					
4. AraProp (Singh, 2022)	0.105	0.600					
5. iCompass (Taboubi et al., 2022)	0.191	0.597					
6. SI2m & AIOX Labs (Gaanoun and Benelallam, 2022)	0.137	0.585					
7. mostafa-samir*	0.186	0.580					
8. Team SIREN AI (Sharara et al., 2022)	0.153	0.578					
9. ChavanKane (Chavan and Kane, 2022)	0.111	0.565					
10. mhmud.fwzi*	0.087	0.552					
11. TUB (Mohtaj and Möller, 2022)	0.076	0.494					
12. tesla*	0.120	0.355					
13. Baseline (Random)	0.043	0.079					
Subtask 2							
1. Pythoneers (Attieh and Hassan, 2022)		0.396					
2. IITD (Mittal and Nakov, 2022)		0.355					
3. NGU_CNLP (Samir et al., 2022)		0.232					
4. Baseline (Random)		0.013					

Random Forests and K-nearest Neighbor. Eventually, for their final submission, they used AraBERT with stacking-based ensemble (5-fold split). They further explored translation-based data augmentation using the English PTC corpus (Da San Martino et al., 2019b).

The second best system was IITD (Mittal and Nakov, 2022), and they used XLM-R and fine-tuned the model. They also explored data augmentation by translating ad adding the PTC corpus as training, but in their experiments this did not help improve the performance.

The third system was CNLP-NITS-PP (Laskar et al., 2022), and they used the AraBERT Twitterbase model along with data augmentation. Note that all systems outperformed the random baseline.

5.2 Subtask 2

In Table 3, we also present an overview of the systems that took part in Subtask 2. Once again, this subtask was dominated by transformer models. We can see in the table that transformers were quite popular, and among them, the most commonly used one was AraBERT, followed by BERT and MARBERT. The participants in this task also used data augmentation and standard pre-processing.

Table 2 shows the evaluation results: we report our random baseline, which is based on the random selection of spans with random lengths and a random assignment of labels.

Table 3: Overview of the approaches used for subtasks 1 and 2, for the teams that submitted a description paper. The systems are ordered by the official score: F1-micro

Rank/Team		Models				Other		
	BERT	XML-R	AraBERT	ARBERT	MARBERT	Data augmentation	Preprocessing	NER
Subtask 1								
1. NGU_CNLP (Samir et al., 2022)	1							_
2. IITD (Mittal and Nakov, 2022)								
3. CNLP-NITS-PP (Laskar et al., 2022)	İ		lacksquare				abla	
3. AraBEM (Eshrag Ali et al., 2022)								
3. Pythoneers (Attieh and Hassan, 2022)			lacksquare				☑	
4. AraProp (Singh, 2022)					abla			
5. iCompass (Taboubi et al., 2022)					✓	_	\square	
6. SI2m & AIOX Labs (Gaanoun and Benelallam, 2022)				abla				
8. Team SIREN AI (Sharara et al., 2022)						_		
9. ChavanKane (Chavan and Kane, 2022)			abla		☑			
11. TUB (Mohtaj and Möller, 2022)	✓						\square	
Subtask 2								
1. Pythoneers (Attieh and Hassan, 2022)			2					_
2. IITD (Mittal and Nakov, 2022)								
3. NGU_CNLP (Samir et al., 2022)								

The best system for this subtask was Pythoneers (Attieh and Hassan, 2022). They used AraBERT with a Conditional Random Field (CRF) layer, which was trained on encoded data using the BIO schema.

The second-best system was IITD (Mittal and Nakov, 2022), which used a Multi-Granularity Network (Da San Martino et al., 2019b) with the mBERT encoder.

The third system was NGU_CNLP (Samir et al., 2022). They converted the data to BIO format and fine-tuned a token classifier based on Marefa-NER³ (pretrained using XLM-RoBERTa).

5.3 Participants' Systems

NGU_CNLP (Samir et al., 2022)[subtask 1:1, subtask 2:3] team participated in both subtasks. For subtask 1, they used a combination of a data augmentation strategy with a transformer-based model. This model ranked first among the 14 systems that participated in this subtask. Their preliminary experiments for subtask 1 consist of using a bag-of-words model with different classical algorithms such as Support Vector Machines, Naïve Bayes, Stochastic Gradient Descent, Logistic regression, Random Forests, and simple K-nearest Neighbor. For subtask 2, they fine-tuned the Marefa-NER model, which is based on XLM-RoBERTa. The system ranked third among the three systems that partici-

³https://huggingface.co/marefa-nlp/marefa-ner

pated in this subtask.

Pythoneers (Attieh and Hassan, 2022)[subtask 1:3, subtask 2:1] also participated in both subtasks. For subtask 1, they trained a multi-task learning model that performs binary classification per propaganda technique. For subtask 2, they first converted the data into BIO format and then fine-tuned an AraBERT model with a Conditional Random Field (CRF) layer. Their subtask 1 system ranked third with a micro-averaged F1-Score of 0.602, and their subtask 2 system ranked first with a micro-averaged F1-Score of 0.396.

IITD (Mittal and Nakov, 2022)[subtask 1:2, subtask 2:2] . This team also participated in both subtasks. They used multilingual pretrained language models for both subtask s. For subtask 1, they used a pretrained XLM-R to estimate a Multinoulli distribution after projecting the CLS embedding to a 20-dimensional embedding (one per propaganda technique). For subtask 2, they used a multigranularity network (Da San Martino et al., 2019b) with mBERT encoder. Even though both systems were trained on only the dataset released in this shared task, they also discussed several methods (zero-shot transfer, continued training, and translation of PTC (Da San Martino et al., 2019b) to Arabic) to study cross-lingual propaganda detection. This suggested interesting research challenges for future exploration, such as how to effectively use data from different domains and how to learn language-agnostic embeddings in propaganda detection systems.

CNLP-NITS-PP (Laskar et al., 2022)[subtask 1:3]. This team participated in subtask 1 and they used AraBERT Twitter-base model for multilabel propaganda classification. They further used data augmentation; in particular, they generated synthetic training data using root and stem substitution from the original train samples and prepared additional synthetic examples. They changed the input labels to the model to be one-hot encoded to indicate multiple labels and modified the macro-F1 scorer to give a score for multiple labels. To make predictions with the model, they used a sentiment analysis pipeline from HuggingFace Transformers and selected all the labels that yielded a score greater than or equal to 0.32. They observed the scores for the predictions on the validation test set and found that most correct labels had a score greater than 0.30. They also found that there was a large gap in the

score for the label when the score was below 0.30.

AraBEM (Eshrag Ali et al., 2022)[subtask 1:3]. This team participated in subtask 1 and they fine-tuned BERT to perform multi-class binary classification. They used standard pre-processing including normalization (mapping letters with various forms, i.e., alef, hamza, and yaa to their representative characters), and removing special characters, diacritics, and repeated characters.

AraProp (Singh, 2022)[subtask 1:4]. This team participated in subtask 1. First, they tokenized the input and produced contextualized word embeddings for all input tokens. To get a fixed-size output representation, they simply averaged all contextualized word embeddings by taking attention mask into account for correct averaging. Then, they added a dropout layer with a dropout rate of 0.3, followed by a linear layer with a sigmoid activation function for the output. They experimented with multiple transformer-based language models: two multilingual models and six monolingual (Arabic) models. Their findings suggest that the MARBERTv2-based fine-tuned model outperforms other models in terms of F1-micro score.

iCompass (Taboubi et al., 2022)[subtask 1:5] team participated in subtask 1. Their system used standard pre-processing such as normalization and removing stopwords, emojis, special characters, and links. Then, they used pre-trained language models such as MARBERT and ARBERT. They further added global average and max pooling layers on top of the models. Finally, they used crossvalidation to improve the model performance.

SI2M & AIOX Labs (Gaanoun and Benelallam, 2022)[subtask 1:6] team participated in subtask 1. They used data augmentation, named entity recognition (NER), and manual rules. For data augmentation, they combined the training and the dev sets, and randomly mixed the sequences to create new synthetic sequences, which they concatenated with the train and the dev sets. Their final system uses a mixed dataset of 2,000 examples. Next, they finetuned ARBERT on the augmented dataset, and they made predictions based on a defined threshold of the classifier's confidence. If no technique got a prediction probability greater than the threshold, the token was assigned the label No technique. Moreover, to detect the Name Calling/Labelling technique, they used a NER model based on AraBERT. Finally, to detect *Repetition*, they used manual rules, after removing the stopwords.

Team SIREN AI (Sharara et al., 2022)[subtask 1:8] participated in subtask 1 and used AraBERT for fine-tuning. Like other teams, they used standard pre-processing, e.g., removing HTML markup, diacritics, non-digit repetitions, etc.

ChavanKane (Chavan and Kane, 2022)[subtask 1:9] team participated in subtask 1 and experimented with AraBERT v1, v02 and v2, MARBERT, ARBERT, XLMRoBERTa, and AraELECTRA. They used a specific variant of DeHateBERT, which is initialized from multilingual BERT and fine-tuned only on Arabic datasets. They also tried creating an ensemble of all models, which consists of five models such as DeHateBERT, AraBERTv2, AraBERTv02, AraBERTv01, and MARBERT. For the final prediction from the ensembles, they used hard voting.

TUB (Mohtaj and Möller, 2022)[subtask 1:11]. This team participated in subtask 1 and used a semantic similarly detection approach based on conceptual word embedding. They converted all sentences in the train, dev, and test sets into vectors using the BERT model. For each sentence in the test set, they detected the five most similar instances from the train and the dev sets, with a cosine similarity above 0.4. Then, they assigned the three most frequent labels among the five instances as the label of the target sentence.

6 Conclusion and Future Work

We presented the WANLP'2022 shared task on Propaganda Detection in Arabic, as part of which we developed the first dataset for Arabic propaganda detection with focus on social media content. This was a successful task: a total of 63 teams registered to participate, and 14 and 3 teams eventually made an official submission on the test set for subtasks 1 and 2, respectively. Finally, 11 teams submitted a task description paper. Subtask 1 asked to identify the propaganda techniques used in a tweet, and subtask 2 further asked to identify the the span(s) of text in which each propaganda technique appears. For both subtasks, the majority of the systems fine-tuned pre-trained Arabic language models, and used standard pre-processing. Some systems used data augmentation and ensemble methods.

In future work, we plan to increase the data size and to add hierarchically structured propaganda techniques.

7 Acknowledgments

This publication was made possible by NPRP grant 13S-0206-200281 Resources and Applications for Detecting and Classifying Polarized and Hate Speech in Arabic Social Media from the Qatar National Research Fund.

Part of this work was also funded by Qatar Foundation's IDKT Fund TDF 03-1209-210013: *Tanbih: Get to Know What You Are Reading.*

This research is also carried out as part of the Tanbih mega-project,⁴ developed at the Qatar Computing Research Institute, HBKU, which aims to limit the impact of "fake news", propaganda, and media bias, thus promoting digital literacy and critical thinking.

The findings herein are solely the responsibility of the authors.

References

Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. 2022. A survey on multimodal disinformation detection. In *Proceedings of the 29th International Conference on Computational Linguistics*, COLING '22, pages 6625–6643, Gyeongju, Republic of Korea.

Firoj Alam, Shaden Shaar, Fahim Dalvi, Hassan Sajjad, Alex Nikolov, Hamdy Mubarak, Giovanni Da San Martino, Ahmed Abdelali, Nadir Durrani, Kareem Darwish, Abdulaziz Al-Homaid, Wajdi Zaghouani, Tommaso Caselli, Gijs Danoe, Friso Stolk, Britt Bruntink, and Preslav Nakov. 2021. Fighting the COVID-19 infodemic: Modeling the perspective of journalists, fact-checkers, social media platforms, policy makers, and the society. In *Findings of EMNLP 2021*, pages 611–649.

Joseph Attieh and Fadi Hassan. 2022. Pythoneers at WANLP 2022 shared task: Monolingual AraBERT for Arabic propaganda detection and span extraction. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, Abu Dhabi, UAE.

Giorgio Barnabò, Federico Siciliano, Carlos Castillo, Stefano Leonardi, Preslav Nakov, Giovanni Da San Martino, and Fabrizio Silvestri. 2022. FbMultiLingMisinfo: Challenging large-scale multilingual benchmark for misinformation detection. In *Proceedings of the 2022 International Joint*

⁴http://tanbih.qcri.org/

- Conference on Neural Networks, IJCNN '22', pages 1–8, Padova, Italy.
- Alberto Barrón-Cedeño, Tamer Elsayed, Preslav Nakov, Giovanni Da San Martino, Maram Hasanain, Reem Suwaileh, Fatima Haouari, Nikolay Babulkov, Bayan Hamdan, Alex Nikolov, Shaden Shaar, and Zien Sheikh Ali. 2020. Overview of CheckThat! 2020 automatic identification and verification of claims in social media. In *Proceedings of the 11th International Conference of the CLEF Association: Experimental IR Meets Multilinguality, Multimodality, and Interaction*, CLEF '2020, pages 215–236, Thessaloniki, Greece.
- Alberto Barrón-Cedeno, 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.
- Tanmay Chavan and Aditya Kane. 2022. Chavankane at WANLP 2022 shared task: Large language models for multi-label propaganda detection. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Giovanni Da San Martino, Alberto Barrón-Cedeño, and Preslav Nakov. 2019a. Findings of the NLP4IF-2019 shared task on fine-grained propaganda detection. In Proceedings of the Second Workshop on Natural Language Processing for Internet Freedom: Censorship, Disinformation, and Propaganda, NLP4IF '19, pages 162–170, Hong Kong, China.
- Giovanni Da San Martino, Alberto Barrón-Cedeno, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020a. SemEval-2020 task 11: Detection of propaganda techniques in news articles. In *Proceedings of the fourteenth workshop on semantic evaluation*, SemEval '20, pages 1377–1414.
- Giovanni Da San Martino, Stefano Cresci, Alberto Barrón-Cedeño, Seunghak Yu, Roberto Di Pietro, and Preslav Nakov. 2020b. A survey on computational propaganda detection. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, IJCAI '20, pages 4826–4832.
- Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, Rostislav Petrov, and Preslav Nakov. 2019b. Fine-grained analysis of propaganda in news article. 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 '19, pages 5636–5646, Hong Kong, China.
- Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021a. Detecting propaganda techniques in memes. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th In-*

- ternational Joint Conference on Natural Language Processing, ACL-IJCNLP '21, pages 6603–6617.
- Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021b. SemEval-2021 task 6: Detection of persuasion techniques in texts and images. In *Proceedings of the 15th International Workshop on Semantic Evaluation*, SemEval '21, pages 70–98.
- Tamer Elsayed, Preslav Nakov, Alberto Barrón-Cedeño, Maram Hasanain, Reem Suwaileh, Pepa Atanasova, and Giovanni Da San Martino. 2019a. CheckThat! at CLEF 2019: Automatic identification and verification of claims. In *Proceedings of the 41st European Conference on Information Retrieval*, ECIR '19, pages 309–315, Cologne, Germany.
- Tamer Elsayed, Preslav Nakov, Alberto Barrón-Cedeño, Maram Hasanain, Reem Suwaileh, Giovanni Da San Martino, and Pepa Atanasova. 2019b. Overview of the CLEF-2019 CheckThat!: Automatic identification and verification of claims. In Experimental IR Meets Multilinguality, Multimodality, and Interaction, LNCS, pages 301–321.
- Refaee Eshrag Ali, Ahmed Basem, and Saad Motaz. 2022. AraBEM at WANLP 2022 shared task: Propaganda detection in Arabic tweets. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Kamel Gaanoun and Benelallam. 2022. SI2M & AIOX Labs at WANLP 2022 shared task: Propaganda detection in Arabic, a data augmentation and name entity recognition approach. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. Transactions of the Association for Computational Linguistics, 10:178–206.
- Ivan Habernal, Raffael Hannemann, Christian Pollak, Christopher Klamm, Patrick Pauli, and Iryna Gurevych. 2017. Argotario: Computational argumentation meets serious games. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 7–12, Copenhagen, Denmark.
- Ivan Habernal, Patrick Pauli, and Iryna Gurevych. 2018. Adapting serious game for fallacious argumentation to German: Pitfalls, insights, and best practices. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, LREC '18, pages 3329–3335, Miyazaki, Japan.
- Andreas Hanselowski, Avinesh PVS, Benjamin Schiller, Felix Caspelherr, Debanjan Chaudhuri, Christian M. Meyer, and Iryna Gurevych. 2018. A retrospective analysis of the fake news challenge

- stance-detection task. In *Proceedings of the 27th International Conference on Computational Linguistics*, COLING '18, pages 1859–1874, Santa Fe, New Mexico, USA.
- Momchil Hardalov, Arnav Arora, Preslav Nakov, and Isabelle Augenstein. 2022. A survey on stance detection for mis- and disinformation identification. In *Findings of the Association for Computational Linguistics: NAACL 2022*, NAACL '22, pages 1259–1277, Seattle, Washington, USA.
- Maram Hasanain, Fatima Haouari, Reem Suwaileh, Zien Sheikh Ali, Bayan Hamdan, Tamer Elsayed, Alberto Barrón-Cedeño, Giovanni Da San Martino, and Preslav Nakov. 2020. Overview of CheckThat! 2020 Arabic: Automatic identification and verification of claims in social media. In Working Notes of CLEF 2020—Conference and Labs of the Evaluation Forum, CLEF '2020, Thessaloniki, Greece.
- Kristina Hristakieva, Stefano Cresci, Giovanni Da San Martino, Mauro Conti, and Preslav Nakov. 2022. The spread of propaganda by coordinated communities on social media. In *Proceedings of the 14th ACM Web Science Conference*, WebSci '22, pages 191–201, Barcelona, Spain.
- Sahinur Rahman Laskar, Rahul Singh, Abdullah Faiz Ur Rahman Khilji, Riyanka Manna, Partha Pakray, and Sivaji Bandyopadhyay. 2022. CNLP-NITS-PP at WANLP 2022 shared task: Propaganda detection in Arabic using data augmentation and AraBERT pre-trained model. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Clyde R. Miller. 1939. The Techniques of Propaganda. From "How to Detect and Analyze Propaganda," an address given at Town Hall. The Center for learning.
- Shubham Mittal and Preslav Nakov. 2022. Iitd at WANLP 2022 shared task: Multilingual multigranularity network for propaganda detection. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Salar Mohtaj and Sebastian Möller. 2022. TUB at WANLP 2022 shared task: Using semantic similarity for propaganda detection in Arabic. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Preslav Nakov, Firoj Alam, Shaden Shaar, Giovanni Da San Martino, and Yifan Zhang. 2021a. COVID-19 in Bulgarian social media: Factuality, harmfulness, propaganda, and framing. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*, RANLP '21, pages 1001–1013, Online. INCOMA Ltd.
- Preslav Nakov, Firoj Alam, Shaden Shaar, Giovanni Da San Martino, and Yifan Zhang. 2021b. A second

- pandemic? Analysis of fake news about COVID-19 vaccines in Qatar. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*, RANLP '21, pages 1014–1025, Online. INCOMA Ltd.
- Preslav Nakov, Alberto Barrón-Cedeño, Giovanni Da San Martino, Firoj Alam, Rubén Míguez, Tommaso Caselli, Mucahid Kutlu, Wajdi Zaghouani, Chengkai Li, Shaden Shaar, Hamdy Mubarak, Alex Nikolov, Yavuz Selim Kartal, and Javier Beltrán. 2022a. Overview of the CLEF-2022 Check-That! lab task 1 on identifying relevant claims in tweets. In Working Notes of CLEF 2022—Conference and Labs of the Evaluation Forum, CLEF '2022, Bologna, Italy.
- Preslav Nakov, Alberto Barrón-Cedeño, Giovanni Da San Martino, Firoj Alam, Julia Maria Struß, Thomas Mandl, Rubén Míguez, Tommaso Caselli, Mucahid Kutlu, Wajdi Zaghouani, Chengkai Li, Shaden Shaar, Gautam Kishore Shahi, Hamdy Mubarak, Alex Nikolov, Nikolay Babulkov, Yavuz Selim Kartal, and Javier Beltrán. 2022b. The CLEF-2022 CheckThat! lab on fighting the COVID-19 infodemic and fake news detection. In Proceedings of the 44th European Conference on IR Research: Advances in Information Retrieval, ECIR '22, pages 416–428, Berlin, Heidelberg.
- Preslav Nakov, Alberto Barrón-Cedeño, Giovanni Da San Martino, Firoj Alam, Julia Maria Struß, Thomas Mandl, Rubén Míguez, Tommaso Caselli, Mucahid Kutlu, Wajdi Zaghouani, Chengkai Li, Shaden Shaar, Gautam Kishore Shahi, Hamdy Alex Nikolov, Nikolay Babulkov, Mubarak, Yavuz Selim Kartal, Javier Beltrán, Michael Wiegand, Melanie Siegel, and Juliane Köhler. 2022c. Overview of the CLEF-2022 CheckThat! lab on fighting the COVID-19 infodemic and fake news detection. In Proceedings of the 13th International Conference of the CLEF Association: Information Access Evaluation meets Multilinguality, Multimodality, and Visualization, CLEF '2022, Bologna, Italy.
- Preslav Nakov, Giovanni Da San Martino, Firoj Alam, Shaden Shaar, Hamdy Mubarak, and Nikolay Babulkov. 2022d. Overview of the CLEF-2022 Check-That! lab task 2 on detecting previously fact-checked claims. In Working Notes of CLEF 2022—Conference and Labs of the Evaluation Forum, CLEF '2022, Bologna, Italy.
- Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeño, Rubén Míguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Nikolay Babulkov, Alex Nikolov, Gautam Kishore Shahi, Julia Maria Struß, and Thomas Mandl. 2021c. The CLEF-2021 CheckThat! lab on detecting check-worthy claims, previously fact-checked claims, and fake news. In *Proceedings of the 43rd European Conference on Information Retrieval*, ECIR '21, pages 639–649, Lucca, Italy.

- Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeño, Rubén Míguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Watheq Mansour, Bayan Hamdan, Zien Sheikh Ali, Nikolay Babulkov, Alex Nikolov, Gautam Kishore Shahi, Julia Maria Struß, Thomas Mandl, Mucahid Kutlu, and Yavuz Selim Kartal. 2021d. Overview of the CLEF-2021 CheckThat! lab on detecting check-worthy claims, previously fact-checked claims, and fake news. In *Proceedings of the 12th International Conference of the CLEF Association: Information Access Evaluation Meets Multiliguality, Multimodality, and Visualization*, CLEF '2021, Bucharest, Romania (online).
- Van-Hoang Nguyen, Kazunari Sugiyama, Preslav Nakov, and Min-Yen Kan. 2022. FANG: Leveraging social context for fake news detection using graph representation. *Commun. ACM*, 65(4):124–132.
- Ray Oshikawa, Jing Qian, and William Yang Wang. 2020. A survey on natural language processing for fake news detection. In *Proceedings of the 12th Language Resources and Evaluation Conference*, LREC '20, pages 6086–6093.
- Konstantin Pogorelov, Daniel Thilo Schroeder, Luk Burchard, Johannes Moe, Stefan Brenner, Petra Filkukova, and Johannes Langguth. 2020. Fake-News: Corona virus and 5G conspiracy task at MediaEval 2020. In *Proceedings of the MediaEval 2020 Workshop*, MediaEval '20.
- Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021. Detecting harmful memes and their targets. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2783–2796.
- 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 *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, EMNLP '17, pages 2931–2937, Copenhagen, Denmark.
- Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Marcos Zampieri, and Preslav Nakov. 2021. SOLID: A large-scale semi-supervised dataset for offensive language identification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP* 2021, pages 915–928.
- Ahmed Samir, Abo Bakr Soliman, Mohamed Ibrahim, Laila Hesham, and Samhaa ElBeltag. 2022. NGU_CNLP at WANLP 2022 shared task: Propaganda detection in Arabic. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Shaden Shaar, Firoj Alam, Giovanni Da San Martino, Alex Nikolov, Wajdi Zaghouani, Preslav Nakov, and

- Anna Feldman. 2021a. Findings of the NLP4IF-2021 shared tasks on fighting the COVID-19 infodemic and censorship detection. In *Proceedings of the Fourth Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda*, NLP4IF '21, pages 82–92.
- Shaden Shaar, Maram Hasanain, Bayan Hamdan, Zien Sheikh Ali, Fatima Haouari, Alex Nikolov, Mucahid Kutlu, Yavuz Selim Kartal, Firoj Alam, Giovanni Da San Martino, Alberto Barrón-Cedeño, Rubén Míguez, Javier Beltrán, Tamer Elsayed, and Preslav Nakov. 2021b. Overview of the CLEF-2021 CheckThat! lab task 1 on check-worthiness estimation in tweets and political debates. Bucharest, Romania (online).
- Shaden Shaar, Alex Nikolov, Nikolay Babulkov, Firoj Alam, Alberto Barrón-Cedeño, Tamer Elsayed, Maram Hasanain, Reem Suwaileh, Fatima Haouari, Giovanni Da San Martino, and Preslav Nakov. 2020. Overview of CheckThat! 2020 English: Automatic identification and verification of claims in social media. In Working Notes of CLEF 2020—Conference and Labs of the Evaluation Forum, CLEF '2020, Thessaloniki, Greece.
- Mohamad Sharara, Wissam Mohamad, Ralph Tawil, Ralph Chobok, Wolf Assi, and Antonio Tannoury. 2022. Team SIREN AI at WANLP 2022 shared task: AraBERT model for propaganda detection. In *Proceedings of the Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Shivam Sharma, Firoj Alam, Md. Shad Akhtar, Dimitar Dimitrov, Giovanni Da San Martino, Hamed Firooz, Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and Tanmoy Chakraborty. 2022. Detecting and understanding harmful memes: A survey. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*, IJCAI '22, pages 5597–5606, Vienna, Austria.
- Gaurav Singh. 2022. AraProp at WANLP 2022 shared task: Leveraging pre-trained language models for Arabic propaganda detection. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- Bilel Taboubi, Bechir Brahem, and Hatem Haddad. 2022. iCompass at WANLP 2022 shared task: ARBERT & MARBERT for multilabel propaganda classification in Arabic tweets. In *Proceedings of the Seventh Arabic Natural Language Processing Workshop*, WANLP '22, Abu Dhabi, UAE.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, NAACL '18, pages 809–819, New Orleans, Louisiana, USA.

Samuel C Woolley and Philip N Howard. 2018. Computational propaganda: political parties, politicians, and political manipulation on social media. Oxford University Press.

Seunghak Yu, Giovanni Da San Martino, Mitra Mohtarami, James Glass, and Preslav Nakov. 2021. Interpretable propaganda detection in news articles. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*, RANLP '21, pages 1597–1605.