

Neural Architectures for Fine-Grained Propaganda Detection in News

Pankaj Gupta^{1,2}, Khushbu Saxena¹, Usama Yaseen^{1,2}, Thomas Runkler¹, Hinrich Schütze²

¹Corporate Technology, Machine-Intelligence (MIC-DE), Siemens AG Munich, Germany

²CIS, University of Munich (LMU) Munich, Germany

pankaj.gupta@siemens.com | pankaj.gupta@campus.lmu.de

Abstract

This paper describes our system (MIC-CIS) details and results of participation in the fine-grained propaganda detection shared task 2019. To address the tasks of sentence (SLC) and fragment level (FLC) propaganda detection, we explore different neural architectures (e.g., CNN, LSTM-CRF and BERT) and extract linguistic (e.g., part-of-speech, named entity, readability, sentiment, emotion, etc.), layout and topical features. Specifically, we have designed multi-granularity and multi-tasking neural architectures to jointly perform both the sentence and fragment level propaganda detection. Additionally, we investigate different ensemble schemes such as majority-voting, relax-voting, etc. to boost overall system performance. Compared to the other participating systems, our submissions are ranked 3rd and 4th in FLC and SLC tasks, respectively.

1 Introduction

In the age of information dissemination without quality control, it has enabled malicious users to spread misinformation via social media and aim individual users with propaganda campaigns to achieve political and financial gains as well as advance a specific agenda. Often disinformation is compiled in the two major forms: fake news and propaganda, where they differ in the sense that the propaganda is possibly built upon true information (e.g., biased, loaded language, repetition, etc.).

Prior works (Rashkin et al., 2017; Habernal et al., 2017; Barrón-Cedeño et al., 2019) in detecting propaganda have focused primarily at document level, typically labeling all articles from a propagandistic news outlet as propaganda and thus, often non-propagandistic articles from the outlet are mislabeled. To this end, Da San Martino et al. (2019) focuses on analyzing the use of propaganda and detecting specific propagandistic

techniques in news articles at sentence and fragment level, respectively and thus, promotes explainable AI. For instance, the following text is a propaganda of type ‘slogan’.

Trump tweeted: “BUILD THE WALL!”
slogan

Shared Task: This work addresses the two tasks in propaganda detection (Da San Martino et al., 2019) of different granularities: (1) *Sentence-level Classification* (SLC), a binary classification that predicts whether a sentence contains at least one propaganda technique, and (2) *Fragment-level Classification* (FLC), a token-level (multi-label) classification that identifies both the spans and the type of propaganda technique(s).

Contributions: (1) To address SLC, we design an ensemble of different classifiers based on Logistic Regression, CNN and BERT, and leverage transfer learning benefits using the pre-trained embeddings/models from FastText and BERT. We also employed different features such as linguistic (sentiment, readability, emotion, part-of-speech and named entity tags, etc.), layout, topics, etc. (2) To address FLC, we design a multi-task neural sequence tagger based on LSTM-CRF and linguistic features to jointly detect propagandistic fragments and its type. Moreover, we investigate performing FLC and SLC jointly in a multi-granularity network based on LSTM-CRF and BERT. (3) Our system (MIC-CIS) is ranked 3rd (out of 12 participants) and 4th (out of 25 participants) in FLC and SLC tasks, respectively.

2 System Description

2.1 Linguistic, Layout and Topical Features

Some of the propaganda techniques (Da San Martino et al., 2019) involve word and phrases that express strong emotional implications, exaggeration, minimization, doubt, national feeling, label-

Category	Feature	Description
<i>Linguistic</i>	POS	part-of-speech tags using NLTK toolkit
	NER	named-entity tags using spacy toolkit, selected tags: {PERSON, NORP, FAC, ORG, GPE, LOC, EVENT, WORK_OF_ART, LAW, LANGUAGE}
	character analysis	count of question and exclamation marks in sentence capital features for each word: first-char-capital, all-char-capital, etc.
	readability	readability and complexity scores using measures from textstat API
	multi-meaning	sum of meanings of a word (grouped by POS) or its synonym nest in the sentence using WordNet
	sentiment	polarity (positive, negative, neutral, compound) scores using spacy; subjectivity using TextBlob; max_pos: maximum of positive, max_neg: max of negative scores of each word in the sentence
	emotional	Emotion features (sadness, joy, fear, disgust, and anger) using IBM Watson NLU API
	loaded words	list of specific words and phrases with strong emotional implications (positive or negative)
<i>Layout</i>	sentence position	categorized as [FIRST, TOP, MIDDLE, BOTTOM, LAST], where, FIRST: 1 st , TOP: < 30%, Middle: between 30-70%, BOTTOM: > 70%, LAST: last sentence of document
	sentence length (l)	categorized as $[= 2, 2 < l \leq 4, 4 < l \leq 8, 8 < l \leq 20, 20 < l \leq 40, 40 < l \leq 60, l > 60]$
<i>Topical</i>	topics	document-topic proportion using LDA, features derived using dominant topic (DT): [DT of current sentence == DT of document, DT of current sentence == DT of the next and previous sentences]
<i>Representation</i>	word vector	pre-trained word vectors from FastText (<i>FastTextWordEmb</i>) and BERT (<i>BERTWordEmb</i>)
	sentence vector	summing word vectors of the sentence to obtain <i>FastTextSentEmb</i> and <i>BERTSentEmb</i>
<i>Decision</i>	relax-boundary	(binary classification) Relax decision boundary and tag propaganda if prediction probability $\geq \tau$
<i>Ensemble</i>	majority-voting	Propaganda if majority says propaganda. In conflict, take prediction of the model with highest F1
	relax-voting	Propaganda if $\mathcal{M} \in [20\%, 30\%, 40\%]$ of models in the ensemble says propaganda.

Table 1: Features used in SLC and FLC tasks

ing, stereotyping, etc. This inspires¹ us in extracting different features (Table 1) including the complexity of text, sentiment, emotion, lexical (POS, NER, etc.), layout, etc. To further investigate, we use topical features (e.g., document-topic proportion) (Blei et al., 2003; Gupta et al., 2019a, 2018) at sentence and document levels in order to determine irrelevant themes, if introduced to the issue being discussed (e.g., *Red Herring*).

For word and sentence representations, we use pre-trained vectors from FastText (Bojanowski et al., 2017) and BERT (Devlin et al., 2019).

2.2 Sentence-level Propaganda Detection

Figure 1 (left) describes the three components of our system for SLC task: features, classifiers and ensemble. The arrows from features-to-classifier indicate that we investigate linguistic, layout and topical features in the two binary classifiers: LogisticRegression and CNN. For CNN, we follow the architecture of Kim (2014) for sentence-level classification, initializing the word vectors by FastText or BERT. We concatenate features in the last hidden layer before classification.

One of our strong classifiers includes BERT that has achieved state-of-the-art performance on mul-

tiple NLP benchmarks. Following Devlin et al. (2019), we fine-tune BERT for binary classification, initializing with a pre-trained model (i.e., *BERT-base, Cased*). Additionally, we apply a decision function such that a sentence is tagged as propaganda if prediction probability of the classifier is greater than a threshold (τ). We relax the binary decision boundary to boost recall, similar to Gupta et al. (2019b).

Ensemble of Logistic Regression, CNN and BERT: In the final component, we collect predictions (i.e., propaganda label) for each sentence from the three ($\mathcal{M} = 3$) classifiers and thus, obtain \mathcal{M} number of predictions for each sentence. We explore two ensemble strategies (Table 1): majority-voting and relax-voting to boost precision and recall, respectively.

2.3 Fragment-level Propaganda Detection

Figure 1 (right) describes our system for FLC task, where we design sequence taggers (Vu et al., 2016; Gupta et al., 2016) in three modes: (1) *LSTM-CRF* (Lample et al., 2016) with word embeddings (w_e) and character embeddings c_e , token-level features (t_f) such as polarity, POS, NER, etc. (2) *LSTM-CRF+Multi-grain* that jointly performs FLC and SLC with FastTextWordEmb and BERTSentEmb, respectively. Here, we add binary

¹some features from datasciencesociety.net/detecting-propaganda-on-sentence-level/

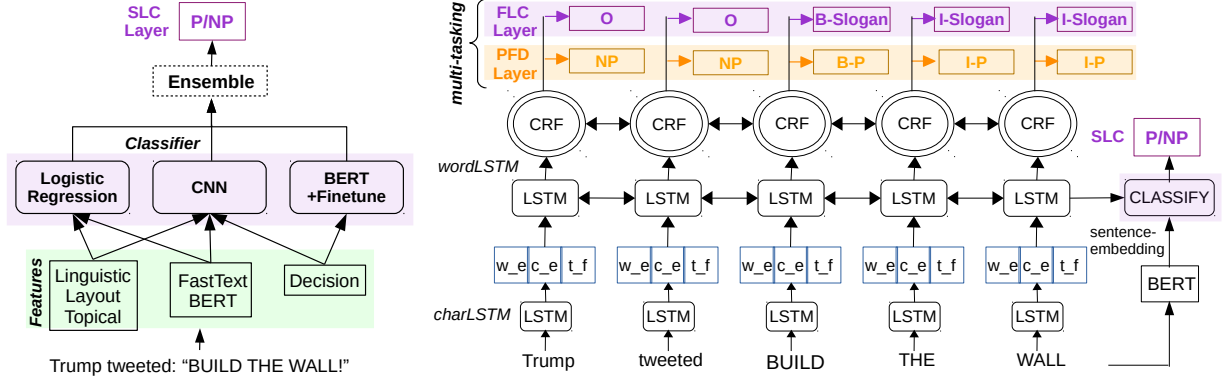


Figure 1: (Left): System description for SLC, including features, transfer learning using pre-trained word embeddings from FastText and BERT and classifiers: LogisticRegression, CNN and BERT fine-tuning. (Right): System description for FLC, including multi-tasking LSTM-CRF architecture consisting of Propaganda Fragment Detection (PFD) and FLC layers. Observe, a binary classification component at the last hidden layer in the recurrent architecture that jointly performs PFD, FLC and SLC tasks (i.e., multi-grained propaganda detection). Here, P: Propaganda, NP: Non-propaganda, B/I/O: Begin, Intermediate and Other tags of BIO tagging scheme.

sentence classification loss to sequence tagging weighted by a factor of α . (3) *LSTM-CRF+Multi-task* that performs propagandistic span/fragment detection (PFD) and FLC (fragment detection + 19-way classification).

Ensemble of Multi-grain, Multi-task LSTM-CRF with BERT: Here, we build an ensemble by considering propagandistic fragments (and its type) from each of the sequence taggers. In doing so, we first perform majority voting at the fragment level for the fragment where their spans exactly overlap. In case of non-overlapping fragments, we consider all. However, when the spans overlap (though with the same label), we consider the fragment with the largest span.

3 Experiments and Evaluation

Data: While the SLC task is binary, the FLC consists of 18 propaganda techniques (Da San Martino et al., 2019). We split (80-20%) the annotated corpus into 5-folds and 3-folds for SLC and FLC tasks, respectively. The development set of each the folds is represented by dev (internal); however, the un-annotated corpus used in leaderboard comparisons by dev (external). We remove empty and single token sentences after tokenization.

Experimental Setup: We use PyTorch framework for the pre-trained BERT model (*Bert-base-cased*²), fine-tuned for SLC task. In the multi-granularity loss, we set $\alpha = 0.1$ for sentence classification based on dev (internal, fold1) scores. We

²github.com/ThilinaRajapakse/pytorch-transformers-classification

Task: SLC (25 participants)			Task: FLC (12 participants)		
Team	F1	P / R	Team	F1	P / R
<i>ltuorp</i>	.6323	.6028 / .6649	<i>newspeak</i>	.2488	.2863 / .2201
<i>ProperGander</i>	.6256	.5649 / .7009	<i>Antiganda</i>	.2267	.2882 / .1869
<i>YMJA</i>	.6249	.6253 / .6246	MIC-CIS	.1999	.2234 / .1808
MIC-CIS	.6231	.5736 / .6819	<i>Stalin</i>	.1453	.1921 / .1169
<i>TeamOne</i>	.6183	.5779 / .6649	<i>TeamOne</i>	.1311	.3235 / .0822

Table 2: Comparison of our system (MIC-CIS) with top-5 participants: Scores on Test set for SLC and FLC

use BIO tagging scheme of NER in FLC task. For CNN, we follow Kim (2014) with filter-sizes of [2, 3, 4, 5, 6], 128 filters and 16 batch-size. We compute binary-F1 and macro-F1³ (Tsai et al., 2006) in SLC and FLC, respectively on dev (internal).

3.1 Results: Sentence-Level Propaganda

Table 3 shows the scores on dev (internal and external) for SLC task. Observe that the pre-trained embeddings (FastText or BERT) outperform TF-IDF vector representation. In row r2, we apply logistic regression classifier with *BERTSentEmb* that leads to improved scores over *FastTextSentEmb*. Subsequently, we augment the sentence vector with additional features that improves F1 on dev (external), however not dev (internal). Next, we initialize CNN by *FastTextWordEmb* or *BERTWordEmb* and augment the last hidden layer (before classification) with *BERTSentEmb* and feature vectors, leading to gains in F1 for both the dev sets. Further, we fine-tune BERT and apply different thresholds in relaxing the decision boundary, where $\tau \geq 0.35$ is found optimal.

³evaluation measure with strict boundary detection

	Dev (internal), Fold1		Dev (external)
	Features	F1 / P / R	F1 / P / R
r1	<i>logisticReg</i> + TF-IDF	.569 / .542 / .598	.506 / .529 / .486
r2	<i>logisticReg</i> + FastTextSentEmb	.606 / .544 / .685	.614 / .595 / .635
	+ Linguistic	.605 / .553 / .667	.613 / .593 / .633
	+ Layout	.600 / .550 / .661	.611 / .591 / .633
	+ Topical	.603 / .552 / .664	.612 / .592 / .633
r3	<i>logisticReg</i> + BERTSentEmb	.614 / .560 / .679	.636 / .638 / .635
r4	+ Linguistic, Layout, Topical	.611 / .564 / .666	.643 / .641 / .644
r5	CNN + FastTextWordEmb	.616 / .685 / .559	.563 / .655 / .494
r6	+ BERTSentEmb	.612 / .693 / .548	.568 / .673 / .491
r7	+ Linguistic, Layout, Topical	.648 / .630 / .668	.632 / .644 / .621
r8	CNN + BERTWordEmb	.610 / .688 / .549	.544 / .667 / .459
r9	+ Linguistic, Layout, Topical	.616 / .671 / .570	.555 / .662 / .478
r10	BERT + Fine-tune ($\tau \geq .50$)	.662 / .635 / .692	.639 / .653 / .625
r11	BERT + Fine-tune ($\tau \geq .40$)	.664 / .625 / .708	.649 / .651 / .647
r12	BERT + Fine-tune ($\tau \geq .35$)	.662 / .615 / .715	.650 / .647 / .654
Ensemble of (r3, r6, r12) within Fold1			
r15	majority-voting $ \mathcal{M} > 50\%$.666 / .663 / .671	.638 / .674 / .605
r16	relax-voting, $ \mathcal{M} \geq 30\%$.645 / .528 / .826	.676 / .592 / .788
Ensemble+ of (r3, r6, r12) from each Fold1-5, i.e., $ \mathcal{M} = 15$			
r17	majority-voting $ \mathcal{M} > 50\%$.666 / .683 / .649
r18	relax-voting, $ \mathcal{M} \geq 40\%$.670 / .646 / .696
r19	relax-voting, $ \mathcal{M} \geq 30\%$		<u>.673</u> / .619 / <u>.737</u>
r20	+ postprocess ($w=10, \lambda \geq .99$)		.669 / .612 / .737
r21	+ postprocess ($w=10, \lambda \geq .95$)		.671 / .612 / .741
Ensemble of (r4, r7, r12) within Fold1			
r22	majority-voting $ \mathcal{M} > 50\%$.669 / .641 / .699	.660 / .663 / .656
r23	relax-voting, $ \mathcal{M} \geq 30\%$.650 / .525 / .852	.674 / .584 / .797
Ensemble+ of (r4, r7, r12) from each Fold1-5, i.e., $ \mathcal{M} = 15$			
r24	majority-voting $ \mathcal{M} > 50\%$.658 / .671 / .645
r25	relax-voting, $ \mathcal{M} \geq 40\%$.673 / .644 / .705
r26	relax-voting, $ \mathcal{M} \geq 30\%$.679 / .622 / .747
r27	+ postprocess ($w=10, \lambda \geq .99$)		.674 / .615 / .747
r28	+ postprocess ($w=10, \lambda \geq .95$)		.676 / .615 / .751

Table 3: SLC: Scores on Dev (internal) of Fold1 and Dev (external) using different classifiers and features.

We choose the three different models in the ensemble: Logistic Regression, CNN and BERT on fold1 and subsequently an ensemble+ of r3, r6 and r12 from each fold1-5 (i.e., 15 models) to obtain predictions for dev (external). We investigate different ensemble schemes (r17-r19), where we observe that the relax-voting improves recall and therefore, the higher F1 (i.e., 0.673). In *postprocess* step, we check for *repetition* propaganda technique by computing cosine similarity between the current sentence and its preceding $w = 10$ sentence vectors (i.e., BERTSentEmb) in the document. If the cosine-similarity is greater than $\lambda \in \{.99, .95\}$, then the current sentence is labeled as propaganda due to repetition. Comparing r19 and r21, we observe a gain in recall, however an overall decrease in F1 applying *postprocess*.

	Dev (internal), Fold1		Dev (external)
	Features	F1 / P / R	F1 / P / R
(I)	<i>LSTM-CRF</i> + FastTextWordEmb	.153 / .228 / .115	.122 / .248 / .081
(II)	+ Polarity, POS, NER	.158 / .292 / .102	.101 / .286 / .061
(III)	+ Multi-grain (SLC+FLC)	.148 / .215 / .112	.119 / .200 / .085
(IV)	+ BERTSentEmb	.152 / .264 / .106	.099 / .248 / .062
(V)	+ Multi-task (PFD)	.144 / .187 / .117	.114 / .179 / .083
Ensemble of (II and IV) within Fold1			
+ postprocess			.116 / .221 / .076
Ensemble of (II and IV) within Fold2			
+ postprocess			.129 / .261 / .085
Ensemble of (II and IV) within Fold3			
+ postprocess			.133 / .220 / .095
Ensemble+ of (II and IV) from each Fold1-3 , i.e., $ \mathcal{M} = 6$			
+ postprocess			.164 / .182 / .150

Table 4: FLC: Scores on Dev (internal) of Fold1 and Dev (external) with different models, features and ensembles. PFD: Propaganda Fragment Detection.

Finally, we use the configuration of r19 on the test set. The ensemble+ of (r4, r7 r12) was analyzed after test submission. Table 2 (SLC) shows that our submission is ranked at 4th position.

3.2 Results: Fragment-Level Propaganda

Table 4 shows the scores on dev (internal and external) for FLC task. Observe that the features (i.e., polarity, POS and NER in row II) when introduced in LSTM-CRF improves F1. We run multi-grained LSTM-CRF without BERTSentEmb (i.e., row III) and with it (i.e., row IV), where the latter improves scores on dev (internal), however not on dev (external). Finally, we perform multi-tasking with another auxiliary task of PFD. Given the scores on dev (internal and external) using different configurations (rows I-V), it is difficult to infer the optimal configuration. Thus, we choose the two best configurations (II and IV) on dev (internal) set and build an ensemble+ of predictions (discussed in section 2.3), leading to a boost in recall and thus an improved F1 on dev (external).

Finally, we use the ensemble+ of (II and IV) from each of the folds 1-3, i.e., $|\mathcal{M}| = 6$ models to obtain predictions on test. Table 2 (FLC) shows that our submission is ranked at 3rd position.

4 Conclusion and Future Work

Our system (Team: MIC-CIS) explores different neural architectures (CNN, BERT and LSTM-CRF) with linguistic, layout and topical features to address the tasks of fine-grained propaganda detection. We have demonstrated gains in per-

formance due to the features, ensemble schemes, multi-tasking and multi-granularity architectures. Compared to the other participating systems, our submissions are ranked 3rd and 4th in FLC and SLC tasks, respectively.

In future, we would like to enrich BERT models with linguistic, layout and topical features during their fine-tuning. Further, we would also be interested in understanding and analyzing the neural network learning, i.e., extracting salient fragments (or key-phrases) in the sentence that generate propaganda, similar to [Gupta and Schütze \(2018\)](#) in order to promote explainable AI.

References

- Alberto Barrón-Cedeño, Giovanni Da San Martino, Israa Jaradat, and Preslav Nakov. 2019. Proppy: A system to unmask propaganda in online news. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019.*, pages 9847–9848.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *TACL*, 5:135–146.
- Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, Rostislav Petrov, and Preslav Nakov. 2019. Fine-grained analysis of propaganda in news articles. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, EMNLP-IJCNLP 2019, Hong Kong, China.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Pankaj Gupta, Yatin Chaudhary, Florian Buettner, and Hinrich Schütze. 2019a. Document informed neural autoregressive topic models with distributional prior. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019.*, pages 6505–6512.
- Pankaj Gupta, Subburam Rajaram, Hinrich Schütze, and Bernt Andrassy. 2018. Deep temporal-recurrent-replicated-softmax for topical trends over time. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 1079–1089.
- Pankaj Gupta, Subburam Rajaram, Hinrich Schütze, and Thomas A. Runkler. 2019b. Neural relation extraction within and across sentence boundaries. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019.*, pages 6513–6520.
- Pankaj Gupta and Hinrich Schütze. 2018. LISA: explaining recurrent neural network judgments via layer-wise semantic accumulation and example to pattern transformation. In *Proceedings of the Workshop: Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2018, Brussels, Belgium, November 1, 2018*, pages 154–164.
- Pankaj Gupta, Hinrich Schütze, and Bernt Andrassy. 2016. Table filling multi-task recurrent neural network for joint entity and relation extraction. In *COLING 2016, 26th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, December 11-16, 2016, Osaka, Japan*, pages 2537–2547.
- Ivan Habernal, Raffael Hannemann, Christian Polak, 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, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017 - System Demonstrations*, pages 7–12.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1746–1751.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural architectures for named entity recognition. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for*

Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016, pages 260–270.

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 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 2931–2937.

Richard Tzong-Han Tsai, Shih-Hung Wu, Wen-Chi Chou, Yu-Chun Lin, Ding He, Jieh Hsiang, Ting-Yi Sung, and Wen-Lian Hsu. 2006. Various criteria in the evaluation of biomedical named entity recognition. *BMC Bioinformatics*, 7:92.

Ngoc Thang Vu, Heike Adel, Pankaj Gupta, and Hinrich Schütze. 2016. Combining recurrent and convolutional neural networks for relation classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 534–539. Association for Computational Linguistics.