

A Retrospective Analysis of the Fake News Challenge Stance Detection Task

Andreas Hanselowski[†], Avinesh PVS[†], Benjamin Schiller[†], Felix Caspelherr[†],
Debanjan Chaudhuri^{‡,*}, Christian M. Meyer[†], Iryna Gurevych[†]

Research Training Group AIPHES

[†]Computer Science Department, Technische Universität Darmstadt

[‡]Smart Data Analytics, University of Bonn

<https://www.aiphes.tu-darmstadt.de>

Abstract

The 2017 Fake News Challenge Stage 1 (FNC-1) shared task addressed a stance classification task as a crucial first step towards detecting fake news. To date, there is no in-depth analysis paper to critically discuss FNC-1’s experimental setup, reproduce the results, and draw conclusions for next-generation stance classification methods. In this paper, we provide such an in-depth analysis for the three top-performing systems. We first find that FNC-1’s proposed evaluation metric favors the majority class, which can be easily classified, and thus overestimates the true discriminative power of the methods. Therefore, we propose a new F1-based metric yielding a changed system ranking. Next, we compare the features and architectures used, which leads to a novel feature-rich stacked LSTM model that performs on par with the best systems, but is superior in predicting minority classes. To understand the methods’ ability to generalize, we derive a new dataset and perform both in-domain and cross-domain experiments. Our qualitative and quantitative study helps interpreting the original FNC-1 scores and understand which features help improving performance and why. Our new dataset and all source code used during the reproduction study are publicly available for future research¹.

1 Introduction

Recently, Pomerleau and Rao (2017) organized the first Fake News Challenge² (FNC-1) in order to foster the development of AI technology to automatically detect fake news. The challenge received much attention in the NLP community: 50 teams from both academia and industry participated. The goal of the FNC-1 challenge is to determine the perspective (or *stance*) of a news article relative to a given headline. An article’s stance can either *agree* or *disagree* with the headline, *discuss* the same topic, or it is completely *unrelated*. Table 1 shows four example documents illustrating these classes.

Stance detection is a crucial building block for a variety of tasks, such as analyzing online debates (Walker et al., 2012; Sridhar et al., 2015; Somasundaran and Wiebe, 2010), determining the veracity of rumors on twitter (Lukasik et al., 2016; Derczynski et al., 2017), or understanding the argumentative structure of persuasive essays (Stab and Gurevych, 2017). While stance detection has been previously focused on individual sentences or phrases, the systems participating in FNC-1 have to detect the stance of an entire document, which raises many new challenges. Although the disagreeing article of Table 1 clearly leans against the headline’s claim, the fourth sentence would agree to it if considered in isolation.

To properly learn from a scientific shared task, there are typically overview and analysis papers that compare the architectures, features, and results of the participating systems. To date, there is, however, no such paper for FNC-1, which is why we conduct a reproduction study of the top three participating systems. Our goal is to independently verify the results reported in the challenge, which is an important asset in empirical research, to critically assess the experimental setup of FNC-1, and to learn building

*The work by Debanjan Chaudhuri has been carried out during his internship at the Ubiquitous Knowledge Processing (UKP) Lab / Adaptive Preparation of Information from Heterogeneous Sources (AIPHES) from 01.01.2017 to 30.06.2017

¹https://github.com/UKPLab/coling2018_fake-news-challenge

²<http://www.fakenewschallenge.org/>

This work is licensed under a Creative Commons Attribution 4.0 International License. License details: <http://creativecommons.org/licenses/by/4.0/>

Headline: Hundreds of Palestinians flee floods in Gaza as Israel opens dams	
Agree (AGR)	GAZA CITY (Ma'an) – Hundreds of Palestinians were evacuated from their homes Sunday morning after Israeli authorities opened a number of dams near the border, flooding the Gaza Valley in the wake of a recent severe winter storm. The Gaza Ministry of Interior said in a statement that civil defense services and teams from the Ministry of Public Works had evacuated more than 80 families from both sides of the Gaza Valley (Wadi Gaza) after their homes flooded as water levels reached more than three meters [...]
Discuss (DSC)	Palestinian officials say hundreds of Gazans were forced to evacuate after Israel opened the gates of several dams on the border with the Gaza Strip, and flooded at least 80 households. Israel has denied the claim as "entirely false". [...]
Disagree (DSG)	Israel has rejected allegations by government officials in the Gaza strip that authorities were responsible for released storm waters flooding parts of the besieged area. "The claim is entirely false, and southern Israel does not have any dams," said a statement from the Coordinator of Government Activities in the Territories (COGAT). "Due to the recent rain, streams were flooded throughout the region with no connection to actions taken by the State of Israel." At least 80 Palestinian families have been evacuated after water levels in the Gaza Valley (Wadi Gaza) rose to almost three meters. [...]
Unrelated (UNR)	Apple is continuing to experience 'Hairgate' problems but they may just be a publicity stunt [...]

Table 1: Headline and text snippets from document bodies with respective stances from the FNC dataset

better methods by understanding their merits and drawbacks. Based on our analysis of the shared task data, we first propose a new evaluation metric for FNC-1 and related document-level stance detection tasks, which is less affected by highly imbalanced datasets. To understand the headroom of the state-of-the-art performance, we additionally estimate the upper bound for this task. In a feature ablation study, we then identify which features contribute to solving the stance detection task. On the basis of our analysis, we combine ideas from previous systems and propose a novel architecture that performs on par with the state-of-the-art systems, but is better able to correctly classify difficult cases. Since generalizability is crucial for the method's future impact, we finally introduce a new evaluation dataset and evaluate how well the FNC-1 models generalize to unseen data from a different domain. In addition to in-domain experiments, we also conduct cross-domain experiments in order to analyze the transfer potential of a method.

2 Related Work

Previous works in stance detection mostly considered target-specific stance prediction, whereby the stance of a text entity with respect to a topic or a named entity is determined. Target-specific stance prediction has been performed for tweets (Mohammad et al., 2016; Augenstein et al., 2016; Zarrella and Marsh, 2016) and online debates (Walker et al., 2012; Somasundaran and Wiebe, 2010; Sridhar et al., 2015). Such target-specific approaches are based on structural (Walker et al., 2012), linguistic and lexical features (Somasundaran and Wiebe, 2010) and they jointly model disagreement only and collective stance using probabilistic soft logic (Sridhar et al., 2015) or neural models (Zarrella and Marsh, 2016; Du et al., 2017) with conditional encoding (Augenstein et al., 2016). **Stance prediction in tweets (Mohammad et al., 2016; Augenstein et al., 2016; Du et al., 2017) and in online debates (Hasan and Ng, 2013) is different from that of stance detection in a news article, which – while similar – is concerned with stance detection of a news article relative to a statement in natural language.**

To the best of our knowledge, there is yet no overview or analysis paper on FNC-1 similar to the shared task on detecting stance in twitter (Mohammad et al., 2016; Derczynski et al., 2017; Taulé et al., 2017). To demonstrate the best scientific practices and achieve research transparency, we close this gap by systematically reviewing the top-ranked systems at FNC-1.

The FNC-1 stance detection task is inspired by Ferreira and Vlachos (2016), who classify the stance of a single sentence of a news headline towards a specific claim. In FNC-1, however, the task is document-level stance detection, which requires the classification of an entire news article relative to a headline. The top performing system in FNC-1 is called *SOLAT in the SWEN* (Sean et al., 2017) by Talos Intelligence (henceforth: Talos). They use a combination of deep convolutional neural networks and gradient-boosted decision trees with lexical features. Team *Athene* (Hanselowski et al., 2017) won the second place with

Dataset	headlines	documents	tokens	instances	AGR	DSG	DSC	UNR
FNC-1	2,587	2,587	372	75,385	7.4%	2.0%	17.7%	72.8%

Table 2: Corpus statistics and label distribution for the FNC-1 dataset

an ensemble of five multi-layer perceptrons (MLP) with six hidden layers each and handcrafted features. For the prediction they used hard voting. Finally, *UCL Machine Reading (UCLMR)* (Riedel et al., 2017) were placed third using a multi-layer perceptron with bag-of-words features. Additionally, recently published work on FNC-1 use a two-step logistic regression based classifier (Bourgonje et al., 2017) and a stacked ensemble of five classifiers (Thorne et al., 2017) which achieve 9th and 11th places respectively. Although multiple systems³ participated at FNC-1, we focus on the top three systems in this paper, due to the availability of source code and our goal of analyzing what contributes most to good performance. In the remaining paper, we introduce and analyze these three systems in detail.

3 Reproduction of the Fake News Challenge FNC-1

In this section, we take a closer look at the challenge. We briefly discuss the task and dataset of FNC-1, describe the three top-ranked systems and reproduce their results.

FNC-1 task and dataset. The task in FNC-1 is learning a classifier $f: (d, h) \mapsto s$ that predicts one of four stance labels $s \in S = \{\text{AGR}, \text{DSG}, \text{DSC}, \text{UNR}\}$ for a document d with regard to a headline h . If headline and document cover different topics, the stance is $s = \text{UNR}$ (unrelated). Otherwise, s is AGR if d agrees and DSG if d disagrees with h . If h and d merely discuss the same topic, but d does not take a definite position, s will be DSC.

To evaluate the challenge, the organizers provide a dataset⁴ of 300 topics. The topics are represented by claims with 5–20 news article documents each. The dataset is derived from the Emergent project (Silverman, 2017) which addressed rumor debunking. In the project, **each news article document was summarized into a headline that reflects the stance of the whole document**. Other than for rumor debunking, the FNC-1 organizers match each document with every summarized headline and then label the (d, h) pair with one of the four stance labels S . To generate the unrelated class UNR, headlines and documents belonging to different topics are randomly matched. Document–headline pairs of 200 topics are reserved for training, the remaining document–headline pairs of 100 topics for testing. Topics, headlines, and documents are therefore not shared between the two data splits. To prevent teams from using any unfair means by deriving the labels for the test set from the publicly available Emergent data, the organizers additionally created 266 instances. Table 2 shows the corpus size and label distribution.

Participating systems. For our reproduction study, we consider FNC-1’s three top-ranked systems. Talos Intelligence’s SOLAT in the SWEN team (Sean et al., 2017) won the FNC-1 using their weighted average model (*TalosComb*) of a deep convolutional neural network (*TalosCNN*) and a gradient-boosted decision trees model (*TalosTree*). *TalosCNN* uses pre-trained word2vec embeddings⁵ passed through several convolutional layers followed by three fully-connected and a final softmax layer for classification. *TalosTree* is based on word count, TF-IDF, sentiment, and singular-value decomposition features in combination with the word2vec embeddings.

Team *Athene* (Hanselowski et al., 2017) was ranked second. They propose a multi-layer perceptron (MLP) inspired by the work of Davis and Proctor (2017). They extend the original model structure to six hidden and a softmax layer and they incorporate multiple hand-engineered features: unigrams, the cosine similarity of word embeddings of nouns and verbs between headline and document tokens, and topic models based on non-negative matrix factorization, latent Dirichlet allocation, and latent semantic indexing in addition to the baseline features provided by the FNC-1 organizers. Depending on the feature type, they either form separate feature vectors for document and headline, or a joint feature vector.

³e.g., <http://web.stanford.edu/class/cs224n/reports.html>

⁴<https://github.com/FakeNewsChallenge/fnc-1-baseline>

⁵<https://code.google.com/archive/p/word2vec/>

Systems	FNC-FNC					
	FNC	F_1m	AGR	DSG	DSC	UNR
Majority vote	.394	.210	0.0	0.0	0.0	.839
TalosComb	.820	.582	.539	.035	.760	.994
TalosTree	.830	.570	.520	.003	.762	.994
TalosCNN	.502	.308	.258	.092	0.0	.882
Athene	.820	.604	.487	.151	.780	.996
UCLMR	.817	.583	.479	.114	.747	.989
featMLP	.825	.607	.530	.151	.766	.982
stackLSTM	.821	.609	.501	.180	.757	.995
Upper bound	.859	.754	.588	.667	.765	.997

Table 3: FNC, F_1m , and class-wise F_1 scores for the analyzed models on in-domain experiments

The UCL Machine Reading (*UCLMR*) team propose an MLP as well, but use only a single hidden layer (Riedel et al., 2017). Their system was ranked third. As features, they use term frequency vectors of unigrams of the 5,000 most frequent words for the headlines and the documents. Additionally, they compute the cosine similarity between the TF-IDF vectors of the headline and document. The resulting term frequency feature vectors of headline and document are concatenated along with the cosine similarity of the two TF-IDF vectors.

Reproduction. Following the instructions from the GitHub repositories of the three teams,⁶ we could successfully reproduce the results reported in the competition without significant deviations. Table 3 shows these results in the FNC column of the *FNC-FNC* setup, which means that the models were trained and tested on the FNC dataset. Since Talos use a combination of two models, we have also included the results of *TalosCNN* and *TalosTree*. A first interesting finding is that *TalosTree* even outperforms the combined model, since the CNN component performed poorly. To understand the merits and drawbacks of the systems, we analyze the performance metrics and the features used, as discussed in the following sections.

4 Performance evaluation

In this section, we critically assess the FNC-1 evaluation methodology and we determine a human upper bound for this task in order to identify the room for improvement for the document-level stance detection task.

Evaluation metrics. The FNC-1 organizers propose the hierarchical evaluation metric FNC, which first awards .25 points if a document is correctly classified as related (i.e., $s \in \{AGR, DSG, DSC\}$) or UNR to a given headline. If it is related, .75 additional points are assigned if the model correctly classifies the document-headline pair as AGR, DSG, or DSC. The goal of this weighting schema is to balance out the large number of unrelated instances.

Nevertheless, the metric fails to take into account the highly imbalanced class distribution of the three related classes AGR, DSG, and DSC illustrated in Table 2. Thus, models, which perform well on the majority class and poorly on the minority classes are favored. Since it is not difficult to separate related from unrelated instances (the best systems reach about $F_1 = .99$ for the UNR class), a classifier that just randomly predicts one of the three related classes would already achieve a high FNC score. A classifier that always predicts DSC for the related documents even reaches $FNC = .833$, which is even higher than the top-ranked system.

We therefore argue that the FNC metric is not appropriate for validating the document-level stance detection task. Instead, we propose the class-wise and the macro-averaged F_1 scores (F_1m) as a new metric for this task that is not affected by the large size of the majority class. The class-wise F_1 scores

⁶<https://github.com/Cisco-Talos/fnc-1>; https://github.com/hanselowski/athene_system; <https://github.com/uclmr/fakenewschallenge>

are the harmonic means of the precisions and recalls of the four classes, which are then averaged to the F_1m metric. The naïve approach of perfectly classifying UNR and always predicting DSC for the related classes, would achieve only $F_1m = .444$, which is clearly different from the proposed systems. By averaging over the individual classes’ scores, F_1m is also applicable to other datasets, which have a different class distribution than the FNC-1 dataset. While the averaged F_1m objectively reflects the quality of the prediction rather than the class distribution, we can also analyze which classes cannot be properly predicted yet.

As the scores in Table 3 indicate, the performance of the three top-ranked systems reach only about $F_1m = .6$. Our analysis reveals that *TalosCNN* does not predict the DSC class yielding an F_1 score of zero. Also the overall performance of this model is low and according to the FNC metric, *TalosTree* would even outperform *TalosComb*. In contrast, *TalosTree* returns almost no predictions for the DSG class, although it performs exceptionally well in terms of FNC. This is because it often predicts the majority class DSC for the related documents. Since there are only few DSG instances in the dataset, the overall performance of this model appears high.

Considering the FNC-1 results according to our proposed F_1m metric, the ranking of the three systems changes: The *TalosComb* and *TalosTree* systems are slightly outperformed by *UCLMR* and clearly outperformed by the Athene system. This is because the two Talos models benefit from the FNC metric definition, favoring the prediction of the majority classes UNR and DSC. On smaller classes, such as DSG, they perform much worse than Athene and *UCLMR*. Using F_1m as a metric, the Athene system would be ranked first, as it outperforms *UCLMR* by 2.1 percentage points. In addition to that, Athene also works best on the DSG, DSC, and UNR class.

Human upper bound. In addition to the issues with the evaluation metric, there is also no upper bound reported for the FNC-1 data, although this will help estimating the headroom of the proposed systems with regard to human performance. Therefore, we ask five human raters to manually label 200 instances. The raters reach an overall inter-annotator agreement of Fleiss’ $\kappa = .686$ (Fleiss, 1971), which is substantial and allows drawing tentative conclusions (Artstein and Poesio, 2008). However, when ignoring the UNR class, the inter-annotator agreement dramatically drops to $\kappa = .218$. This indicates that differentiating between the three related classes AGR, DSG, and DSC is difficult even for humans.

On the basis of the annotation, we also determine the most probable stance labels according to MACE (Hovy et al., 2013), and compare them to the ground truth from the Emergent project. The agreement of the labels in this case is better, reaching a Fleiss’ κ of .807 overall and .552 for the three related classes. The MACE-based most probable label allows us to compute the human upper bound as $F_1m = .754$, which we include in Table 9 along with the upper bound per class F1 scores UNR = .997 AGR = .588, DSG = .667, and DSC = .765.

5 Analysis of models and features

In this section, we first perform an error analysis in order to be able to find out what the three best performing models are learning and in which cases they fail. In order to address the identified drawbacks, we conduct a systematic feature analysis and derive an alternative model based on our findings.

5.1 Error analysis

Our error analysis for the three analyzed systems shows that the models fail in the following cases: (1) If there is lexical overlap between the headline and the document, the models classify the instance as one of the related classes, even in cases in which the two are unrelated. (2) If the document–headline pair is related, but only contains synonyms rather than the same tokens, the model often misclassifies the case as UNR. (3) If keywords like *reports*, *said*, or *allegedly* are detected, the systems classify the pair as DSC. (4) The DSG class is especially difficult to determine, as only few lexical indicators (e.g., *false*, *hoax*, *fake*) are available as features. The disagreement is often expressed in complex terms which demands more sophisticated machine learning techniques. For example: “If the bizarre story about...sounded outlandish, that’s because it was”. In appendix A.2, we illustrate these errors with concrete examples.

The analysis shows that the models exploit the similarity between the headline and the document in terms of lexical overlap. Lexical cue words, such as *reports*, *said*, *false*, *hoax* play an important role in classification. However, the systems fail when semantic relations between words need to be taken into account, complex negation instances are encountered, or the understanding of propositional content in general is required. This is not surprising since the three models are based on n -grams, bag-of-words, topic models and lexicon-based features instead of capturing the semantics of the text. In this section, we test these features systematically and we propose new features and a new architecture for FNC-1.

5.2 Feature analysis

Throughout our feature analysis, we use the Athene model, which performed best in terms of F_1 m and allows a large number of experiments due to its fast computation. All tests are performed on the FNC-1 development set with 10-fold cross-validation. In the remaining section, we first discuss and evaluate the performance of each feature individually and then conduct an ablation test for groups of similar features. Detailed feature descriptions are included in the supplementary material (section A.1). Figure 1 shows the system performance of the individual features discussed below.

FNC-1 baseline features. The FNC-1 organizers provide a gradient-boosting baseline using the co-occurrence (COOC) of word and character n -grams in the headline and the document as well as two lexicon-based features, which count the number of refuting (REFU) and polarity (POLA) words based on small word lists. Figure 1 indicates that COOC performs well, whereas both lexicon-based features are on par with the majority vote baseline.

Challenge features. The three analyzed FNC-1 systems rely on combinations of the following features: Bag-of-words (BoW) unigram features, topic model features based on non-negative matrix factorization (NMF-300, NMF-cos) (Lin, 2007), Latent Dirichlet Allocation (LDA-cos) (Blei et al., 2001), Latent Semantic Indexing (LSI-300) (Deerwester et al., 1990), two lexicon-based features using NRC Hashtag Sentiment (NRC-Lex) and Sentiment140 (Sent140) (Mohammad et al., 2013), and word similarity features which measure the cosine similarity of pre-trained word2vec embeddings of nouns and verbs in the headlines and the documents (WSim). The topic models use 300 topics. Besides the concatenated topic vectors, we also consider the cosine similarity between the topics of document and headline (NMF-cos, LDA-cos). The BoW features perform best in terms of F_1 m. While LSI-300, NMF-300 and NMF-cos topic models yield high scores, LDA-cos and WSim fall behind.

Novel features. We also analyze a number of novel features for the FNC-1 task which have not been used in the challenge. Bag-of-character 3-grams (BoC) represent subword information. They show promising results in our setup. The structural features (STRUC) include the average word lengths of the headline and the document, the number of paragraphs in the document and their average lengths. The low performance of this feature indicates that the structure of the headline and the documents is not indicative of their stance. Furthermore, we test readability features (READ) which estimate the complexity of a text. Less complex texts could be indicative of deficiently written fake news. We tried the following metrics for headline and document as a concatenated feature vector: SMOG grade (Mc Laughlin, 1969), Flesch-Kincaid grade level and Flesch reading ease (Kincaid et al., 1975), Gunning fog index (Štajner et al., 2012), Coleman-Liau index (Mari and Ta Lin, 1975), automated readability index (Senter and Smith, 1967), LIX and RIX (Jonathan, 1983), McAlpine EFLAW Readability Score (McAlpine, 1997), and Strain Index (Solomon, 2006). However, in the present problem setting these features show only a low performance. The same is true for the lexical diversity (LexDiv) metrics, type-token ratio, and the measure of textual diversity (MTLD) (McCarthy, 2005). We finally analyze the performance of features based on the following lexicons: MPQA (Wilson et al., 2005), MaxDiff (Kiritchenko et al., 2014), and EmoLex (Mohammad and Turney, 2010). These features are based on the sentiment, polarity, and emotion expressed by headlines and documents, which might be good indicators of an author’s opinion. However, our results show that these lexicon-based features are not helpful. Even though the considered lexicons are important for fake-news detection (Shu et al. (2017), Horne and Adali (2017)), for stance detection, the properties captured by the lexicon-based features are not very useful.

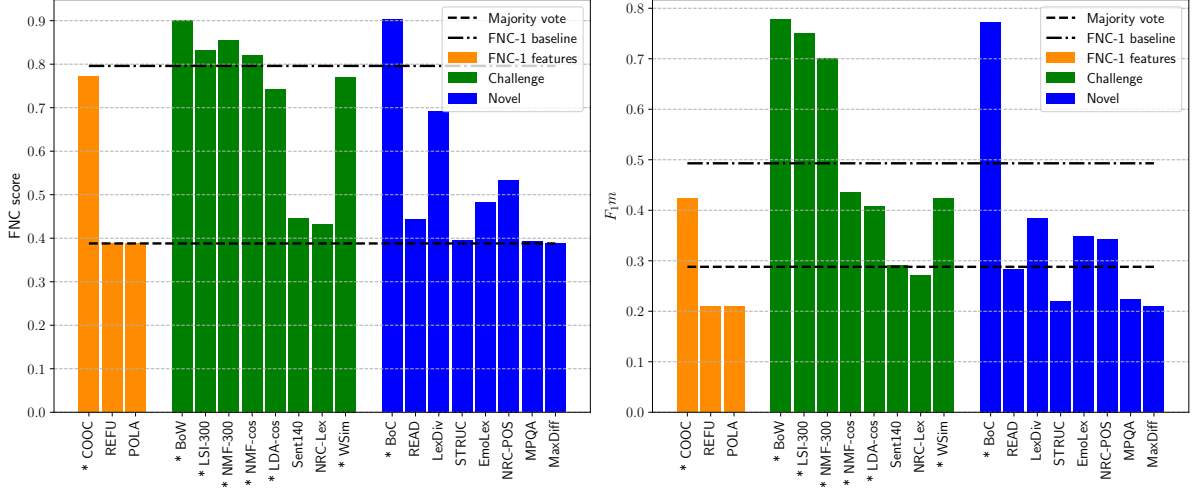


Figure 1: Performance of the system based on individual features

Feature ablation test. We first remove all features that are more than 10% below the FNC-1 baseline, since they mostly predict the majority class and thus harm the F_{1m} score. In Figure 1, we mark these features with an asterisk (*). To quantify their contribution, we perform an ablation test across three groups of related features: (1) BoW and BoC (BoW/C), (2) LSI-topic, NMF-topic, NMF-cos, LDA-cos (Topic), and (3) NRC-POS and WSim (Oth).

Table 4 show the results of our ablation test. The BoW and BoC features have the biggest impact on the performance. While the topic models yield further improvements, the NRC-POS and WSim features are not helpful. Hence, we suggest BoW, BoC, and the four topic model based features as the most promising feature set.

We evaluate this feature set on the FNC-1 test dataset. The results are included in the *featMLP* row of Table 3 for the *FNC-FNC* setting. Although *featMLP* with the revised feature selection outperforms the best performing FNC-1 system Athene in terms of F_{1m} and FNC score, the margin is not significant. Similar to the three FNC-1 systems, we observe a .2 performance drop between the development and test dataset. This is most likely because of the 100 new topics in the test dataset, which have not been seen during training. Thus, the evaluation on the test set can be considered as an out-of-domain prediction.

5.3 Model analysis

In order to increase the overall performance, we conduct additional experiments with an ensemble of the three models *featMLP*, *TalosComb*, and *UCLMR* using hard voting. However, we could not significantly improve the results. Since all models struggle with the DSG class, we have applied different under- and over-sampling techniques to balance the class distribution, but also this technique did not yield improved results.

In the error analysis, we observed that the feature-based systems lack semantic understanding. Therefore, we combine a feature-based system with a model that is better able to capture the semantics based

	Baselines		Only			All without			All*	All
	Maj. vote	FNC-1	BoW/C	Topic	Oth	-BoW/C	-Topic	-Oth		
AGR	0.0	.241	.772	.637	0.0	.665	.714	.722	.713	.675
DSG	0.0	.047	.601	.571	0.0	.530	.598	.616	.573	.455
DSC	0.0	.738	.874	.838	.731	.841	.863	.876	.870	.835
UNR	.835	.970	.991	.983	.964	.982	.989	.995	.993	.989
F_{1m}	.209	.499	.796	.757	.425	.754	.791	.802	.787	.738

Table 4: Results of the feature ablation test. Baseline FNC-1 uses gradient boosting classifier with all FNC-1 baseline features. * states that only the preselected features are used (see Figure 1).

Headline: NHL expansion ahead? No, says Gary Bettman

Article body: It wasn't very long ago that NHL commissioner Gary Bettman was treating talk of expansion as though he was being asked if he'd like an epidemic of Ebola. But recently the nature of the rhetoric has changed so much that the question is becoming not if, but when. ...

Table 5: A correctly classified DSG instance by the *stackLSTM*

on word embeddings and sequential encoding. Sequential processing of information is important in order to get the meaning of the whole sentence, e.g. "It wasn't long ago that Gary Bettman was ready to expand NHL." VS. "It was long ago that Gary Bettman wasn't ready to expand NHL." In Figure 2, we introduce this *stackLSTM* model, which combines the best feature set found in the ablation test with a stacked long short-term memory (LSTM) network (Hermans and Schrauwen, 2013). We use 50-dimensional GloVe word embeddings⁷ (Pennington et al., 2014) in order to generate sequences of word vectors of a headline–document pair. For this, we concatenate a maximum of 100 tokens of the headline and the document. These embedded word sequences v_1, v_2, \dots, v_n are fed through two stacked LSTMs with a hidden state size of 100 with a dropout of 0.2 each. The last hidden state of the second LSTM is concatenated with the feature set and fed into a 3-layer neural network with 600 neurons each. Finally, we add a dense layer with four neurons and softmax activation function in order to retrieve the class probabilities.

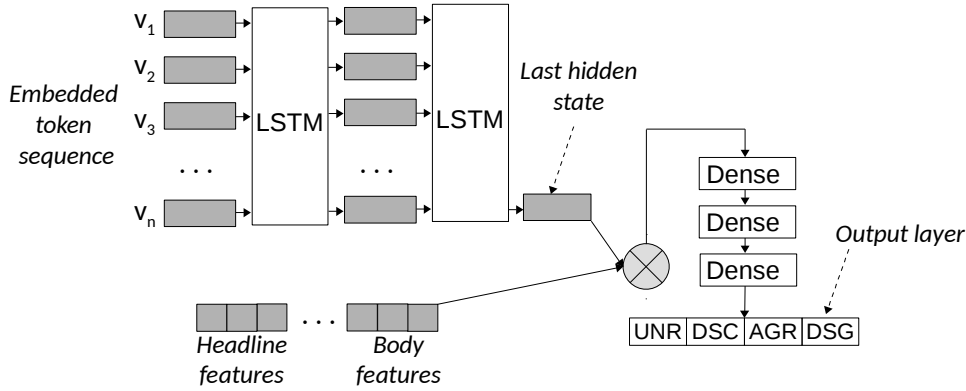


Figure 2: Model Architecture of the feature-rich *stackLSTM*

Table 3 shows the performance of *stackLSTM* for the *FNC-FNC* setup. Our model outperforms all other methods in terms of F_1m . The difference to *Athene* and *featMLP* is, however, not significant. An important advantage of *stackLSTM* is its improved performance for the DSG class, which is the most difficult one to predict due to the low number of instances. This means that *stackLSTM* correctly classifies a larger number of complex negation instances. The difference on this difficult DSG class between *stackLSTM* and all other methods is statistically significant (using Student's t-test). The model predicts more often for the DSG class and gets more of these examples correct without compromising the overall performance. One challenging DSG example, which was correctly classified by the *stackLSTM*, is given in Table 5.

6 Analysis of the generalizability of the models

To test the robustness of the models (i.e. how well they generalize to new datasets), we introduce novel test data for document-level stance detection based on the Argument Reasoning Comprehension (ARC) task proposed by Habernal et al. (2018). In this section, we describe the dataset, analyze the models' performance, and perform cross-domain experiments.

ARC dataset. Habernal et al. (2018) manually select 188 debate topics with popular questions from the user debate section of the New York Times. For each topic, they collect user posts, which are highly ranked by other users, and create two claims representing two opposing views on the topic. Then, they

⁷<http://nlp.stanford.edu/data/glove.twitter.27B.zip>

Dataset	headlines	documents	tokens	instances	AGR	DSG	DSC	UNR
ARC	4,448	4,448	99	17,792	8.9%	10.0%	6.1%	75.0%

Table 6: Corpus statistics and label distribution for the ARC dataset

Example from the original ARC dataset		
Topic	Do same-sex colleges play an important role in education, or are they outdated?	
User post	Only 40 women’s colleges are left in the U.S. And, while there are a variety of opinions on their value, to the women who have attended ... them, they have been ... tremendously valuable. ...	
Claims	1. Same-sex colleges are outdated 2. Same-sex colleges are still relevant	
Label	Same-sex colleges are still relevant	
Generated instance in alignment with the FNC problem setting		
Stance	Headline	Document
AGR	Same-sex colleges are still relevant	Only 40 women’s colleges are left in the U.S. ...

Table 7: Example of the original ARC dataset and the generated instance to align with FNC dataset

ask crowd workers to decide whether a user post supports either of the two opposing claims or does not express a stance at all. This Argument Reasoning Comprehension (ARC) dataset consists of typical controversial topics from the news domain, such as *immigration*, *schooling issues*, or *international affairs*. While this is similar to the FNC-1 dataset, there are significant differences, as a user post is typically a multi-sentence statement representing one viewpoint on the topic. In contrast, the news articles of FNC-1 are longer and usually provide more balanced and detailed perspective on an issue.

To allow using the ARC data for our FNC stance detection setup, we consider each user post as a document and randomly select one of the two claims as the headline. We label the claim–document pair as AGR if the claim has also been chosen by the workers, as DSG if the workers chose the opposite claim, and as DSC if the workers selected neither claim. Table 7 shows an example of our revised ARC corpus structure. In order to generate the unrelated instances, we randomly match the user posts with claims, but avoid that a user post is assigned to a claim from the same topic. Table 6 provides basic corpus statistics. For training and testing, we split the corpus into 80% training/validation set and 20% testing set.

As for the FNC-1 corpus, we have also determined a human upper bound for the ARC dataset. Five subjects annotate 200 samples using the four classes. Even though the overall Fleiss’ $\kappa = .614$ is slightly lower compared to the FNC-1 corpus, the agreement for the three related classes AGR, DSG, and DSC is higher ($\kappa = .383$) than for FNC-1. The human upper bound based on MACE is $F_{1m} = .773$. Table 9 contains also the class-wise F_1 scores.

In-domain experiments ARC-ARC: The in-domain results for the ARC corpus listed in Table 8 show that the overall performance of all models decreases. Since the models have been constructed to perform well on the FNC-1 dataset, this is not surprising. Nevertheless, for the ARC corpus, the models are better able to distinguish between AGR and DSG instances. We assume this is because the number of DSG instances is substantially larger and is similar to the number AGR instances. The classification of the DSC instances, on the other hand, turns out to be more challenging on the ARC corpus. This is because even if a user post is related to the claim, it often does not explicitly refer to it. With *TalosComb* being best, the Talos models were better able to generalize to the new data. Even though the *stackLSTM* is again better on the more difficult minority class (in this case DSC), the structure and features of *TalosComb* seem to be more appropriate for this problem setting.

Cross-domain experiments: In the cross-domain setting we train on the training data of one corpus and evaluate on the test data of the other corpus. The experiments in Table 9 show that the performance of the models is substantially better than the majority vote baseline. We therefore conclude that the two problem settings are related and exhibit a common structure. The results suggest that *TalosComb* is best able to learn from the ARC corpus, as it is also superior in the *ARC-FNC* setting. The *stackLSTM*, on the

Systems	ARC-ARC					
	FNC	F_1m	AGR	DSG	DSC	UNR
Majority vote	.430	.214	0.0	0.0	0.0	.857
TalosComb	.725	.573	.593	.598	.160	.944
Athene	.680	.548	.516	.482	.190	.933
UCLMR	.667	.519	.517	.503	.121	.932
featMLP	.690	.526	.526	.506	.144	.934
stackLSTM	.685	.524	.451	.518	.194	.935
Upper bound	.796	.773	.710	.857	.571	.954

Table 8: FNC, F_1m , and class-wise F_1 scores for the analyzed models on in-domain experiments

Systems	FNC-ARC						ARC-FNC					
	FNC	F_1m	AGR	DSG	DSC	UNR	FNC	F_1m	AGR	DSG	DSC	UNR
Majority vote	.430	.214	0.0	0.0	0.0	.857	.394	.210	0.0	0.0	0.0	.839
TalosComb	.584	.365	.336	0.0	.195	.929	.607	.388	.279	.183	.113	.977
Athene	.523	.340	.340	.244	.138	.894	.548	.321	.277	.097	.028	.882
UCLMR	.557	.358	.271	.064	.201	.896	.482	.288	.234	.109	.080	.728
featMLP	.586	.389	.321	.159	.171	.906	.585	.351	.322	.111	.033	.939
stackLSTM	.591	.401	.321	.191	.182	.910	.613	.373	.343	.116	.082	.950
Upper bound	.796	.773	.710	.857	.571	.954	.859	.754	.588	.667	.765	.997

Table 9: FNC, F_1m and class-wise F_1 scores(F_1m) based on cross-domain experiments

other hand, yields best results when trained on the FNC corpus as the *FNC-ARC* setting suggests.

7 Discussion and conclusion

In this paper, we conducted a thorough analysis of the Fake News Challenge stage one stance detection task. Although this is common for shared tasks, there is yet no analysis or reproduction study of this task, which is why we close this gap. Given that the challenge has attracted much attention in the NLP community with 50 participating teams, a detailed analysis is valuable as it provides insights into the problem setting and lessons learned for upcoming competitions. In our investigation, we evaluated the performance of the three top-scoring systems, critically assessed the experimental setup, and performed a detailed feature analysis, in which we identify high-performing features for the task yielding a new model. We conducted an error analysis and found that the models mostly rely on the lexical overlap for classification. To assess how well the models generalize to a similar problem setting, we experimented with a second, newly derived corpus. We also propose a new evaluation metric based on F_1 scores, since the challenge’s metric is highly affected by the imbalanced class distribution of the test data. Using this evaluation setup, the ranking of the top three systems changes. Based on our analysis, we conclude that the investigated stance detection problem is challenging, since the best performing features are not yet able to resolve the difficult cases. Thus, more sophisticated machine learning techniques are needed, which have a deeper semantic understanding, and are able to determine the stance on the basis of propositional content instead of relying on lexical features.

8 Acknowledgements

This work has been supported by the German Research Foundation as part of the Research Training Group “Adaptive Preparation of Information from Heterogeneous Sources” (AIPHES) at the Technische Universität Darmstadt under grant No. GRK 1994/1.

References

Ron Artstein and Massimo Poesio. 2008. Inter-coder agreement for computational linguistics. *Computational Linguistics* 34(4):555–596.

- Isabelle Augenstein, Tim Rocktäschel, Andreas Vlachos, and Kalina Bontcheva. 2016. Stance detection with bidirectional conditional encoding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Austin, TX, USA, pages 876–885.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2001. Latent dirichlet allocation. In *Advances in Neural Information Processing Systems 14 (NIPS)*. Vancouver, BC, Canada, pages 601–608.
- Peter Bourgonje, Julian Moreno Schneider, and Georg Rehm. 2017. From Clickbait to Fake News Detection: An Approach based on Detecting the Stance of Headlines to Articles. In *Proceedings of the EMNLP 2017 Workshop ‘Natural Language Processing meets Journalism’*. Copenhagen, Denmark, pages 84–89.
- Sanjiv R. Das and Mike Y. Chen. 2007. Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web. *Management Science* 53(9):1375–1388.
- Richard Davis and Chris Proctor. 2017. Fake News, Real Consequences: Recruiting Neural Networks for the Fight Against Fake News. Online: <http://web.stanford.edu/class/cs224n/reports/2761239.pdf>. Accessed: 2018-03-16.
- Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American society for information science* 41(6):391.
- Leon Derczynski, Kalina Bontcheva, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Arkaitz Zubiaga. 2017. SemEval-2017 Task 8: RumourEval: Determining rumour veracity and support for rumours. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval)*. Vancouver, BC, Canada, pages 69–76.
- Jiachen Du, Ruifeng Xu, Yulan He, and Lin Gui. 2017. Stance classification with target-specific neural attention. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI)*. Melbourne, Australia, pages 3988–3994.
- William Ferreira and Andreas Vlachos. 2016. Emergent: a novel data-set for stance classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL/HLT)*. San Diego, CA, USA, pages 1163–1168.
- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76(5):378.
- Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL/HLT)*. New Orleans, LA, USA, pages 1930–1940.
- Andreas Hanselowski, Avinesh PVS, Benjamin Schiller, and Felix Caspelherr. 2017. Description of the system developed by team Athene in the FNC-1, 2017. Online: https://github.com/hanselowski/athene_system/blob/master/system_description_athene.pdf. Accessed: 2018-03-13.
- Kazi Saidul Hasan and Vincent Ng. 2013. Stance Classification of Ideological Debates: Data, Models, Features, and Constraints. In *Proceedings of the Sixth International Joint Conference on Natural Language Processing (IJCNLP)*. Nagoya, Japan, pages 1348–1356.
- Michiel Hermans and Benjamin Schrauwen. 2013. Training and analysing deep recurrent neural networks. In *Advances in neural information processing systems 26 (NIPS)*. Stateline, NV, USA, pages 190–198.
- Benjamin D. Horne and Sibel Adali. 2017. This Just In: Fake News Packs a Lot in Title, Uses Simpler, Repetitive Content in Text Body, More Similar to Satire than Real News. In *Proceedings of the ICWSM 2017 Workshop on News and Public Opinion*. Montréal, QC, Canada, pages 759–766.

- Dirk Hovy, Taylor Berg-Kirkpatrick, Ashish Vaswani, and Eduard Hovy. 2013. Learning Whom to Trust with MACE. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL/HLT)*. Atlanta, GA, USA, pages 1120–1130.
- Anderson Jonathan. 1983. Lix and Rix: Variations on a Little-known Readability Index. *Journal of Reading* 26(6):490–496.
- J. Peter Kincaid, Robert P. Fishburne Jr, Richard L. Rogers, and Brad S. Chissom. 1975. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. Research Branch Report 8-75, Naval Technical Training Command, Millington, TN, USA.
- Svetlana Kiritchenko, Xiaodan Zhu, and Saif M. Mohammad. 2014. Sentiment analysis of short informal texts. *Journal of Artificial Intelligence Research* 50:723–762.
- Chih-Jen Lin. 2007. Projected gradient methods for nonnegative matrix factorization. *Neural Computation* 19(10):2756–2779.
- Michal Lukasik, P.K. Srijith, Duy Vu, Kalina Bontcheva, Arkaitz Zubiaga, and Trevor Cohn. 2016. Hawkes processes for continuous time sequence classification: an application to rumour stance classification in twitter. In *Proceedings of 54th Annual Meeting of the Association for Computational Linguistics (ACL)*. Berlin, Germany, pages 393–398.
- Coleman Mari and Liao Ta Lin. 1975. A computer readability formula designed for machine scoring. *Journal of Applied Psychology* 60(2):283.
- G. Harry Mc Laughlin. 1969. SMOG grading—a new readability formula. *Journal of reading* 12(8):639–646.
- Rachel McAlpine. 1997. *Global English for global business*. Longman.
- Philip M. McCarthy. 2005. An assessment of the range and usefulness of lexical diversity measures and the potential of the measure of textual, lexical diversity (MTLD). *Dissertation Abstracts International* 66:12.
- Saif M. Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. SemEval-2016 Task 6: Detecting Stance in Tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval)*. San Diego, CA, USA, pages 31–41.
- Saif M. Mohammad, Svetlana Kiritchenko, and Xiaodan Zhu. 2013. NRC-Canada: Building the State-of-the-Art in Sentiment Analysis of Tweets. In *Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval)*. Atlanta, GA, USA, pages 321–327.
- Saif M. Mohammad and Peter D. Turney. 2010. Emotions Evoked by Common Words and Phrases: Using Mechanical Turk to Create an Emotion Lexicon. In *Proceedings of the NAACL/HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*. Los Angeles, CA, USA, pages 26–34.
- Saif M. Mohammad and Peter D. Turney. 2013. Crowdsourcing a word-emotion association lexicon 29(3):436–465.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar, pages 1532–1543.
- Dean Pomerleau and Delip Rao. 2017. The Fake News Challenge: Exploring how artificial intelligence technologies could be leveraged to combat fake news. <http://www.fakenewschallenge.org/>. Accessed: 2017-10-20.
- Benjamin Riedel, Isabelle Augenstein, Georgios P Spithourakis, and Sebastian Riedel. 2017. A simple but tough-to-beat baseline for the fake news challenge stance detection task. *arXiv preprint arXiv:1707.03264*.

- Sara Rosenthal, Preslav Nakov, Svetlana Kiritchenko, Saif Mohammad, Alan Ritter, and Veselin Stoyanov. 2015. SemEval-2015 Task 10: Sentiment Analysis in Twitter. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval)*. Denver, CO, USA, pages 451–463.
- Baird Sean, Sibley Doug, and Pan Yuxi. 2017. Talos Targets Disinformation with Fake News Challenge Victory. <http://blog.talosintelligence.com/2017/06/talos-fake-news-challenge.html>. Accessed: 2017-12-02.
- R.J. Senter and Edgar A. Smith. 1967. Automated readability index. Technical Report AMRL-TR-66-220, University of Cincinnati.
- Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter* 19(1):22–36.
- Craig Silverman. 2017. Emergent: A real-time rumor tracker. Online: <http://www.emergent.info/>. Accessed: 2017-12-13.
- N. Watson Solomon. 2006. *Strain Index: A New Readability Formula*. Master thesis, Madurai Kamaraj University.
- Swapna Somasundaran and Janyce Wiebe. 2010. Recognizing stances in ideological on-line debates. In *Proceedings of the NAACL/HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*. Los Angeles, CA, USA, pages 116–124.
- Dhanya Sridhar, James R. Foulds, Bert Huang, Lise Getoor, and Marilyn A. Walker. 2015. Joint models of disagreement and stance in online debate. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing (ACL/IJCNLP)*. Beijing, China, pages 116–125.
- Christian Stab and Iryna Gurevych. 2017. Parsing argumentation structures in persuasive essays. *Computational Linguistics* 43(3):619–659.
- Sanja Štajner, Richard Evans, Constantin Orăsan, and Ruslan Mitkov. 2012. What can readability measures really tell us about text complexity. In *Proceedings of the LREC 2012 Workshop on Natural Language Processing for Improving Textual Accessibility (NLP4ITA)*. Istanbul, Turkey, pages 14–21.
- Mariona Taulé, Maria Antònia Martí, Francisco M. Rangel Pardo, Paolo Rosso, Cristina Bosco, and Viviana Patti. 2017. Overview of the Task on Stance and Gender Detection in Tweets on Catalan Independence. In *Proceedings of the Second Workshop on Evaluation of Human Language Technologies for Iberian Languages (IberEval) co-located with 33th Conference of the Spanish Society for Natural Language Processing (SEPLN)*. Murcia, Spain, pages 157–177.
- James Thorne, Mingjie Chen, Giorgos Myrianthous, Jiashu Pu, Xiaoxuan Wang, and Andreas Vlachos. 2017. Fake news stance detection using stacked ensemble of classifiers. In *Proceedings of the EMNLP 2017 Workshop ‘Natural Language Processing meets Journalism’*. Copenhagen, Denmark, pages 80–83.
- Marilyn A. Walker, Pranav Anand, Robert Abbott, and Ricky Grant. 2012. Stance classification using dialogic properties of persuasion. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL/HLT)*. Montréal, QC, Canada, pages 592–596.
- Theresa Wilson, Janyce Wiebe, and Paul Hoffmann. 2005. Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing (HLT/EMNLP)*. Vancouver, BC, Canada, pages 347–354.
- Guido Zarrella and Amy Marsh. 2016. MITRE at SemEval-2016 Task 6: Transfer Learning for Stance Detection. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. San Diego, CA, USA, pages 458–463.

Xiaodan Zhu, Svetlana Kiritchenko, and Saif M. Mohammad. 2014. NRC-Canada-2014: Recent Improvements in the Sentiment Analysis of Tweets. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval)*. Dublin, Ireland, pages 443–447.

A Supplemental Material

A.1 Features: Detailed description

BoW/BoC features We use bag-of-words (BoW) 1- and 2-grams with 5,000 tokens vocabulary for the headline as well as the document. For the BoW feature, based on a technique by Das and Chen (2007), we add a negation tag "_NEG" as prefix to every word between special negation keywords (e.g. "not", "never", "no") until the next punctuation mark appears. For the bag-of-characters (BoC) 3-grams are chosen with 5,000 tokens vocabulary, too. For the BoW/BoC feature we use the TF to extract the vocabulary and to build the feature vectors of headline and document. The resulting TF vectors of headline and document get concatenated afterwards. Feature *co-occurrence* (FNC-1 baseline feature) counts how many times word 1-/2-/4-grams, character 2-/4-/8-/16-grams, and stop words of the headline appear in the first 100, first 255 characters of the document, and how often they appear in the document overall.

Topic models We use non-negative matrix factorization (NMF) (Lin, 2007), latent semantic indexing (LSI) (Deerwester et al., 1990), and latent Dirichlet allocation (LDA) (Blei et al., 2001) to create topic models out of which we create independent features. For each topic model, we extract 300 topics out of the headline and document texts. Afterwards, we compute the similarity of headlines and bodies to the found topics separately and either concatenate the feature vectors (NMF, LSI) or calculate the cosine distance between them as a single valued feature (NMF, LDA).

Lexicon-based features These features are based on the NRC Hashtag Sentiment and Sentiment140 lexicon (Kiritchenko et al., 2014; Mohammad et al., 2013; Zhu et al., 2014), as well as for the MPQA lexicon (Wilson et al., 2005) and MaxDiff Twitter lexicon (Rosenthal et al., 2015; Kiritchenko et al., 2014). All named lexicons hold values that signal the sentiment/polarity for each word. The features are computed separately for headline and document, and constructed as proposed by Mohammad et al. (2013): First, we count how many words with positive, negative, and without polarity are found in the text. Two features sum up the positive and negative polarity values of the words in the texts and another two features are set by finding the word with the maximum positive and negative polarity value in the text. Finally, the last word in the text with negative or positive polarity is taken as a feature. Since the MaxDiff Twitter lexicon also contains 2-grams, we decide to take them into account as well, whereas for the other lexicons only 1-grams incorporated. Additionally, we base features on the EmoLex lexicon (Mohammad and Turney, 2010, 2013). For all its words, it holds up to eight emotions (anger, fear, anticipation, trust, surprise, sadness, joy, disgust), based on the context they frequently appear in. For headline and document respectively, the emotions for all words are counted as a feature vector. The resulting vectors for headline and document are then concatenated. Lastly, the baseline features *polarity words* and *refuting words* are added. The first one counts refuting words (e.g. "fake", "hoax"), divides the sum by two, and takes the remainder as a feature signaling the polarity of headline or document. The latter one sets a binary feature for each refuting word (e.g. "fraud", "deny") appearing in the headline or document.

Readability features We measure the readability of headline and document with SMOG grade (only document), Flesch-Kincaid grade level, Flesch reading ease, and Gunning fog index (Štajner et al., 2012), Coleman-Liau index (Mari and Ta Lin, 1975), automated readability index (Senter and Smith, 1967), LIX and RIX (Jonathan, 1983), McAlpine EFLAW Readability Score (McAlpine, 1997), Strain Index (Solomon, 2006). The SMOG grade is only valid if a text has at least 30 sentences, and thus is only implemented for the bodies.

Lexical features As lexical features we implement the type-token-ratio (TTR) and the measure of textual lexical diversity (MTLD) (McCarthy, 2005) for the document, and only type-token-ratio for the headline, since MTLD needs at least 50 tokens to be valid. Also, the baseline feature *word overlap* belongs to this group. It divides the cardinality of the intersection of unique words in headline and document by the cardinality of the union of unique words in headline and document.

POS features The POS features amongst others include counters for nouns, personal pronouns, verbs and verbs in past tense, adverbs, nouns and proper nouns, cardinal numbers, punctuations, the ratio of quoted words, and also the frequency of the three least common words in the text. The headline feature also contains a value for the percentage of stop words and the number of verb phrases, which showed good results in the work of Horne and Adali (2017). For the *word-similarity* feature, which are mainly based on Ferreira and Vlachos (2016) we calculated average word embeddings (pre-trained word2vec model⁸) for all verbs (retrieved with Stanford Core NLP toolkit⁹) of headline/document separately. The cosine similarity between the averaged embeddings of headline and document is taken as a feature, as well as the hungarian distance between verbs of headline and document based on the paraphrase database¹⁰. The same computation is done for all nouns of headline and document. Additionally the average sentiment of the headline and the average sentiment of the document is used as a feature. A count of negating words of the headline and the document is added to the feature vector as well as the distance from the negated word to the root of the sentence. The number of average words per sentence of headline and document is another feature. The aforementioned features are improved by only selecting a predefined number of sentences of document and headline. Therefore the sentences are ordered by TF-IDF score.

Structural features The structural features contain the average word length of the headline and document, and the number of paragraphs and average paragraph length of the document.

A.2 Misclassified examples identified in the error analysis

Example 1.

(ground truth: "unrelated", system predicts: "agree")

Headline: CNN: Doctor Took Mid-Surgery Selfie with Unconscious Joan Rivers

Document: "A TEENAGER woke up during brain surgery to ask doctors how it was going. Iga Jasica, 19, was having an op to remove a tumour at when the anaesthetic wore off and she struck up a conversation with the medics still working on her."

Example 2.

(ground truth: "agree", system predicts: "unrelated")

Headline: Three Boobs Are Most Likely Two Boobs and a Lie

Document: The woman who claimed she had a third breast has been proved a hoax.

Example 3.

(ground truth: "disagree", system predicts: "discuss")

Headline: Woman pays 20,000 for third breast to make herself LESS attractive to men

Document: The woman who reported that she added a third breast was most likely lying.

Example 4.

(ground truth: "disagree", system predicts: "agree")

⁸<https://code.google.com/archive/p/word2vec/>

⁹<https://stanfordnlp.github.io/CoreNLP/>

¹⁰<http://www.cis.upenn.edu/~ccb/ppdb/>

Headline: Disgusting! Joan Rivers Doc Gwen Korovin's Sick Selfie EXPOSED — Last Photo Of Comic Icon, When She Was Under Anesthesia

Document: If the bizarre story about Joan Rivers' doctor pausing to take a "selfie" in the operating room minutes before the 81-year-old comedienne went into cardiac arrest on August 29 sounded outlandish, that's because it was.