

POSHAN: Cardinal POS Pattern Guided Attention for News Headline Incongruence

Rahul Mishra

University of Stavanger, Norway
rahul.mishra@uis.no

Shuo Zhang

Bloomberg, United Kingdom
szhang611@bloomberg.net

ABSTRACT

Automatic detection of click-bait and incongruent news headlines is crucial to maintaining the reliability of the Web and has raised much research attention. However, most existing methods perform poorly when news headlines contain contextually important cardinal values, such as a quantity or an amount. In this work, we focus on this particular case and propose a neural attention based solution, which uses a novel cardinal **Part of Speech (POS)** tag pattern based hierarchical attention network, namely **POSHAN**, to learn effective representations of sentences in a news article. In addition, we investigate a novel cardinal phrase guided attention, which uses word embeddings of the contextually-important cardinal value and neighbouring words. In the experiments conducted on two publicly available datasets, we observe that the proposed method gives appropriate significance to cardinal values and outperforms all the baselines. An ablation study of POSHAN shows that the cardinal POS-tag pattern-based hierarchical attention is very effective for the cases in which headlines contain cardinal values.

CCS CONCEPTS

• **Information systems** → **Document representation**; • **Computing methodologies** → **Neural networks**.

KEYWORDS

News Headline Incongruence; Cardinal Part-Of-Speech Pattern; Neural Attention

ACM Reference Format:

Rahul Mishra and Shuo Zhang. 2021. POSHAN: Cardinal POS Pattern Guided Attention for News Headline Incongruence. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM '21)*, November 1–5, 2021, Virtual Event, QLD, Australia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3459637.3482376>

1 INTRODUCTION

News titles expose the first impression to readers and decide the viral potential of news stories within social networks [22]. Most of the users only rely on the news title to decide what to read further [10]. A deceptive and misleading news title can lead to false beliefs and wrong opinions. It becomes inversely worse when users share

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '21, November 1–5, 2021, Virtual Event, QLD, Australia

© 2021 Association for Computing Machinery.
ACM ISBN 978-1-4503-8446-9/21/11...\$15.00
<https://doi.org/10.1145/3459637.3482376>

Headline: *Immigration Expert : US Will Have 100 Million New Immigrants in Next 50 Years.*

Body: *These projections show that new immigrants and their descendants will drive most U.S. population growth in the coming 50 years, as they have for the past half-century. Among the projected 441 million Americans in 2065, 78 million will be immigrants and 81 million will be people born in the U.S. to immigrant parents.*

Figure 1: Example of an Incongruent Headline. The headline says, in the next 50 years, there will be 100 million new immigrants but the news body quotes about only 78 million new immigrants.

the news on social media without reading the news body but only skimming through the news title. The news headlines, which are ambiguous, misleading, and deliberately made catchy to lure the users to click, are called incongruent headlines or click baits [28]. Figure 1 illustrates an example. There is a line of study that has been investigated and analyzed in the literature [4, 5, 9, 20, 23, 26, 28, 30], witnessed by different techniques such as linguistic feature based methods [5], generative adversarial networks [17, 23], and hierarchical neural attention networks [30]. For example, Yoon et al. [30] propose a headline text guided neural attention network to compute an incongruence score between the news headline and the corresponding body text. They introduce a hierarchical attention based encoder, which encodes words of news body text at the word level to form paragraph representations and encodes paragraphs to form document representation. However, we observe that these prior works fail to generalize and perform adequately in cases where news headlines contain a significant numerical value. The numerical values can be in the form of a currency amount, counts of people, months, years or objects, etc. For instance, in Figure 1 an excerpt from a news item is shown, in which the news headline “*Immigration Expert : US Will Have 100 Million New Immigrants in Next 50 Years.*”, contains two contextually important numerical figures i.e. “100 Million”, “50 Years.”. The headline mentions, there will be 100 million new immigrants, but the news body quotes only 78 million new immigrants. The headline is deliberately made contradictory and exaggerating to look more sensational. It’s apparent from this example that numerical and cardinal values are useful and crucial cues of the congruence of the news headlines.

All of the prior works suffer from not giving enough importance to numerical values. The headline guided attention-based methods such as [30], fail to attend relevant words related to cardinal phrases as they do not treat them specifically. On the other hand, generative adversarial network-based methods, which generate a synthetic

headline from news body text to augment the dataset or to use them for similarity matching with original headlines, also miss cardinal aspects in the synthetically generated headlines. Clearly, the news headlines having numbers are not trivial cases for incongruence detection, and in this paper, we try to deal with news headline incongruence detection with special focus to such cases.

The objective of this work is to devise an incongruence detection method, which not only performs better than previously proposed techniques but also resolves the deterioration of classification accuracy with the news items in which the headline contains cardinal values. In specific, we leverage a novel Cardinal Part-of-Speech Tag patterns to drive the hierarchical neural attention to capture salient and contextually important words and sentences at the word and sentence levels correspondingly. The key idea of using cardinal pos patterns such as (*NN : CD : JJ*) or (*VBD : CD : CD*) is to use them as latent features associated with news headlines and learn the contextual embeddings based on data samples containing the same cardinal pos patterns. The embeddings are used at the test time to drive the attention and capture the salient words and sentences, which are significant for cardinal values. In addition, we investigate a cardinal phrase guided attention mechanism and combine both with standard headline guided attention. To utilize the better contextual representation of words, we fine tune the pre-trained BERT model and extract the word embeddings, which are fed to a Bi-LSTM based sequence encoder. We conduct experiments with the subset of two publicly available datasets and achieve state-of-the-art performance. The proposed model *POSHAN* not only outperforms all the other methods in original datasets but also its performance does not deteriorate much compared to other models, with derived datasets, containing only those data samples, which have numerical values in the news headlines. We visualize the Cardinal POS Pattern embeddings and overall attention weights to further analyse the effectiveness of the proposed model. We observe that the Cardinal POS Patterns have formed clearly separated clusters in embedding space, which connote the congruence and incongruence labels. It is apparent from the visualization of overall attention weights, that *POSHAN* model successfully attends the contextually important cardinal phrases in addition to other significant words. In nutshell, the major contributions of this work are:

- We focus on the news headline incongruence detection when news headlines containing numbers, and propose the cardinal POS pattern guided attention (Section 3.4) baseline.
- We propose a cardinal phrase guided attention (Section 3.6) mechanism and combine the both cardinal POS pattern and cardinal phrase attention with standard headline guided attention (Section 3.7) in a joint model (Section 3.8).
- We incorporate the proposed hierarchical attention methods on top of a Bi-LSTM based sequence encoder (Section 3.3) which encodes the sequence of fine-tuned (Section 5.2) pre-trained BERT embeddings (Section 3.2) of the words.
- In the evaluation with two publicly available datasets (Table 3 and 4), the proposed techniques outperform the baselines and state-of-the-art methods.
- We visualize the Cardinal POS Pattern embeddings and overall attention weights and conduct error analysis to analyze

the effectiveness of the proposed model, and verify the effectiveness.

2 RELATED WORK

Detection and prevention of misinformation and deceptive content online has gained lots of traction recently. Incongruent news and click-baits are very common forms of deception and misinformation. Naturally, most of the prior works in this area have treated the click-baits or news incongruence detection task as a standard text classification problem. Majority of the initial works are feature engineering heavy [3], exploiting diverse features such as linguistic features, lexicons, sentiments and statistical features. Chakraborty et al. [4] use linguistic and syntactic features such as sentence structure, word patterns, word n-grams and part-of-speech (POS) n-grams etc. and learn a classifier using support vector machine (SVM) to detect click-baits. Potthast et al. [20] use text features and meta-information of tweets such as entity mentions, emotional polarity, tweet length and word n-grams to learn a classifier, experimenting with methods such as random forest, logistic regression etc. to detect click-baits. Chen et al. [5] propose to conduct lexical and syntactic analysis and advocate to utilize image features and user-behavior features for the identification of the click-baits. These methods are outperformed by the recent deep learning based methods [1, 13], in which hand crafted feature engineering is not required. News headline incongruence is closely related to a number of tasks such as sentence matching based stance classification [9, 26]. Sentence pair classification task using fine tuned pre-trained language models such as BERT [8] and RoBERTa [14] has received a great traction from the community and it is a closely related problem to headline incongruence. Sentence pair classification typically consists a pair of sentences, while in headline incongruence systems, we need to deal with a sentence and a large news body content in order to form the evidence for congruence. The sentence matching and lexical similarity based methods [18] are not a good fit for headline incongruence problem due to inherent challenges such as relative length and vocabulary mismatch between the news headline text and its body content. Therefore, these tasks share the advancement of the development of the techniques. For example, Wei and Wan [28] introduce a co-training based approach with myriad kinds of features such as sentiments, textual, and informality. Recent works such as [30] use neural attention [2] based approach to achieve headline guided contextual representation of the news body text and also release a Korean and an English dataset for headline incongruence.

Although some recent works such as [30] have achieved state-of-the-art performance, most of the existing approaches do not perform well in the case of the headline containing cardinal numbers because no additional emphasis is given to the cardinal numbers. Shu et al. [23] use a generative approach to augment the dataset by additionally generating synthetic headlines. Mishra et al. [17] use an inter-mutual attention-based semantic matching between the original and a synthetically generated headlines via generative adversarial network based techniques, which utilises the difference between all the pairs of word embeddings of words involved and computes mutual attention score matrix. These generative methods are also not very useful in the task at hand as the news

headlines, generated using news body content, usually miss the cardinal information. Focusing on the quantity cases, we propose a neural attention mechanism in which, we use novel cardinal POS triplet and cardinal phrase guided attention in addition to standard headline guided attention. This technique makes sure to have two contextual information: firstly, by using headline guided attention, all the keywords of the headline are utilized in forming the overall attention oriented representation. Secondly, by applying cardinal POS triplet and cardinal phrase guided attention, we ensure that cardinal value is emphasized and overall representation contains the effect of cardinal value.

3 PROBLEM DEFINITION AND PROPOSED MODEL

In this section, we formally introduce the problem definition. Then we present the overall architecture of the proposed model *POSHAN* sequentially, see Figure 2. In specific, we first discuss the embedding layer, which outputs the vector representations of the words and cardinal pos-tag patterns. Secondly, we describe the cardinal pos-tag pattern guided hierarchical attention in detail for both word and sentence level. Thirdly, we introduce cardinal phrase guided and headline text guided hierarchical attention. In the end, we explain a method to fuse all the three attention types to get the overall attention scores.

3.1 Problem Definition

Given a news item $n_i \in N$, where N is the set of all news items, which has a headline h_i and body content b_i , *news title incongruence detection* aims to predict the news as ‘‘Congruent (C)’’ or ‘‘Incongruent (I)’’, where incongruence denotes a mismatch between h_i and b_i by content. News headline h_i and news body content b_i are comprising of sequence of l words as $h_i = \{w_{h1}, w_{h2}, \dots, w_{hl}\} \in W$ and m words as $b_i = \{w_{b1}, w_{b2}, \dots, w_{bm}\} \in W$ correspondingly, where W is the overall vocabulary set.

3.2 Embedding Layer

In Figure 2, for a news item n_i , the corresponding headline h_i of length l and body content b_i of length m are represented as $h_i = \{f(w_{h1}), \dots, f(w_{hl})\}$ where $f(w_{hj}) \in \mathbb{R}^d$ is a word embedding vector of dimension d for j^{th} word in headline h_i and $b_i = \{f(w_{b1}), f(w_{b2}), \dots, f(w_{bm})\}$ where $f(w_{bk}) \in \mathbb{R}^d$ is a word embedding vector of dimension d for k^{th} word in body content b_i . We use pre-trained contextual BERT embeddings, extracted using bert-as-service [29] tool to get the embeddings of the size of 768 dimensions for each word. Each headline is associated with a cardinal pos-tag pattern of form $(POS_p POS_q POS_r)$, where POS_p , POS_q and POS_r are the pos-tags corresponding to the cardinal phrase $(W_p W_q W_r)$. We describe the cardinal pos-tag pattern and cardinal phrase in detail in sections 3.4 and 3.6 respectively. We also create the vector representation \vec{PP} for each of the cardinal pos-tag patterns of the size of 100 dimensions and initialize them with uniformly random weights. We learn weights for these cardinal pos-tag pattern embeddings jointly in the *POSHAN* model via backprop of error, as shown in the Figure 3.

3.3 Sequence Encoder

The pre-trained contextual BERT embeddings of the words of the news body text $b_i = \{f(w_{b1}), f(w_{b2}), \dots, f(w_{bm})\}$ are fed to a Bi-directional Long short term memory (Bi-LSTM) unit [11] based encoder, which encodes the news body text using standard LSTM equations. The output from Bi-LSTM units are the concatenations of forward and backward hidden states for each word, i.e., $HS_{i,j} = \overrightarrow{hs_{i,j}} \parallel \overleftarrow{hs_{i,j}}$. Where $\overrightarrow{hs_{i,j}}$ and $\overleftarrow{hs_{i,j}}$ are the forward and backward hidden states of Bi-LSTM units. $HS_{i,j}$ is the overall hidden state for the i^{th} word of the j^{th} sentence.

3.4 Cardinal POS Triplet Patterns

The idea of utilizing POS-tag Patterns to capture the intended context in natural language text is inspired by prior works [12, 19]. Justeson and Katz [12] propose and utilize 7 handcrafted part-of-speech (POS) patterns to extract significant and useful phrases from a long unstructured text. We utilize part-of-speech patterns containing cardinal POS tag ‘CD’ and call it cardinal POS triplet patterns. A cardinal POS triplet pattern can be defined as $(* : CD : *)$, where in place of wildcards, there can be (JJ, NN, VB) etc., e.g. $(NN : CD : JJ)$. In contrast to [12], we do not handcraft a list of the viable POS patterns rather we use the all possible combination of POS patterns of length 3, containing POS tag ‘CD’. We apply a neural attention layer in which, these cardinal POS triplets are used to guide the attention to select salient words and sentences which are significant for the POS pattern.

3.5 Cardinal POS Triplet Pattern Guided Hierarchical Attention

The objective of the Cardinal POS Triplet Pattern attention is to attend or select salient words that are significant and have some connotation with the cardinal phrase of the headline. Similarly, we aim to attend the salient sentences at the sentence level attention. Yoon et al. [30] have used headline guided attention to model the contextual representation of the news body text. However, we observe that the headline guided attention is not sufficient and effective, in case of headlines containing cardinal values. During experiments, we noticed that only headline-based attention convolutes the effective representation and fails to capture the influence of cardinal phrases on the overall document representation. We take a different and more logical design decision, in which we use part-of-speech patterns contained in each headline h_i to guide the attention. We learn an embedding \vec{PP} for each cardinal POS triplet pattern as discussed in section 3.2.

Computing Word Level Attention weights: We use the embedding of cardinal POS triplet pattern \vec{PP} to compute the attention scores given to each hidden state of the Bi-LSTM encoder.

$$S^j = \sum_{i=1}^l \alpha_{ij}^p HS_{ij} \quad (1)$$

Where HS_{ij} is the hidden state for the i^{th} word of j^{th} sentence and l is maximum number of words in a sentence. α_{ij}^p is the attention weight. S^j is the formed sentence representation of j^{th} sentence after attention scores are applied. The attention score α_{ij} can be defined as:

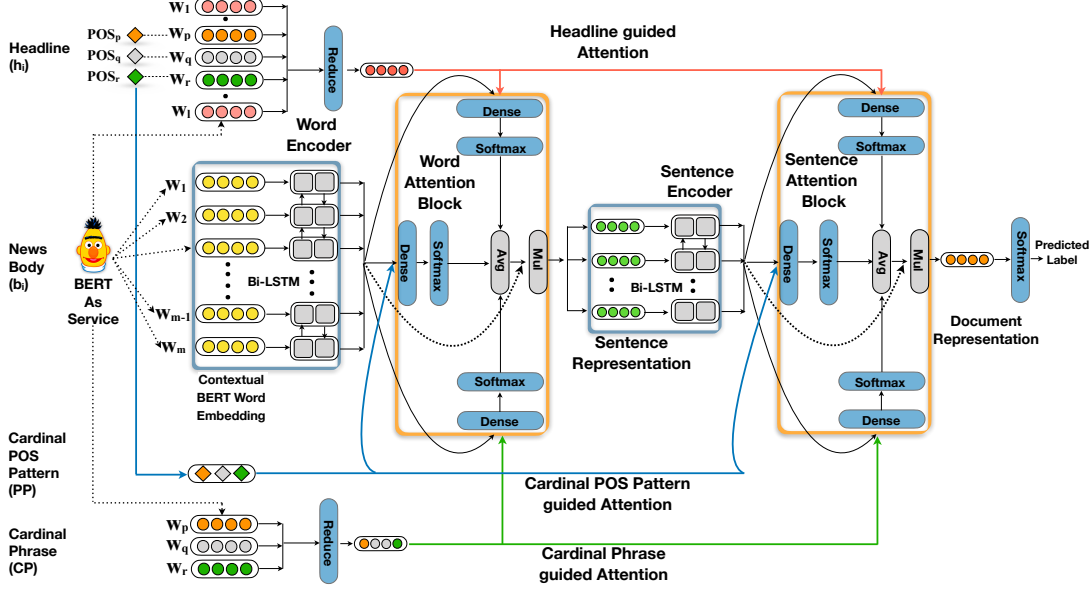


Figure 2: Overall Architecture of POSHAN Model: Rectangle with blue borders are the Bi-LSTM based encoder at the word and sentence levels. Rectangle with orange borders are the Attention blocks at the word and sentence levels. Headline guided Attention, Cardinal POS Pattern guided Attention and Cardinal Phrase guided Attention are depicted as red, blue and green connecting lines respectively.

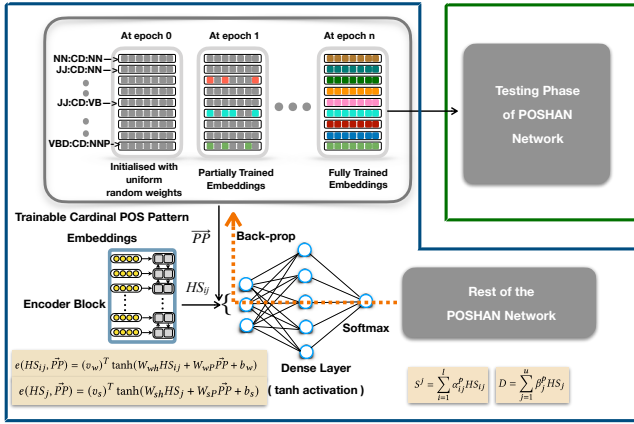


Figure 3: Training of Cardinal POS Pattern Embeddings: Inspired by a recent work [16], we create trainable Cardinal POS Pattern embeddings of 100 dimensions for each POS pattern and initialize them with the uniformly random weights to get the representation of POS patterns in vector space. The weights of these embeddings are trained during training of POSHAN model via back-prop of error. At the test time, we use already trained Cardinal POS Pattern embeddings, trained during training.

$$\alpha_{ij}^p = \frac{\exp(e(HS_{ij}, \vec{PP}))}{\sum_{k=1}^l \exp(e(HS_{ik}, \vec{PP}))} \quad (2)$$

Where e is a \tanh based scoring function, which is used to compute the attention scores. \vec{PP} is the POS-Tag pattern vector. The scoring function $e(HS_{ij}, \vec{PP})$ can be defined as:

$$e(HS_{ij}, \vec{PP}) = (v_w)^T \tanh(W_{wh}HS_{ij} + W_{wp}\vec{PP} + b_w) \quad (3)$$

Where v_w is weight vector at the word level. W_{wh} and W_{wp} are the weight matrices for hidden state and aspect vector and b_w is bias at the word level respectively.

Computing Sentence Level Attention weights: To compute sentence level POS-Tag pattern driven attention weights, we use POS-Tag pattern vector representation \vec{PP} and hidden states HS_j^S from the sentence level BI-LSTM units as concatenations of both forward and backward hidden states $HS_j^S = \overrightarrow{hs}_j^S \parallel \overleftarrow{hs}_j^S$ as follows:

$$D = \sum_{j=1}^o \beta_j^p HS_j \quad (4)$$

Where HS_j is the hidden state for j^{th} sentence and β_j^p is the attention weight. o is the maximum no of sentences in a news body text. D is the formed document representation after the attention scores are applied. The attention score β_j can be defined as:

$$\beta_j^p = \frac{\exp(e(HS_j, \vec{PP}))}{\sum_{k=1}^o \exp(e(HS_k, \vec{PP}))} \quad (5)$$

Where e is a \tanh based scoring function, which is used to compute the attention scores. \vec{PP} is the POS-Tag pattern vector. The scoring function $e(HS_j, \vec{PP})$ can be defined as:

$$e(HS_j, \vec{PP}) = (v_s)^T \tanh(W_{sh}HS_j + W_{sp}\vec{PP} + b_s) \quad (6)$$

Where v_s is weight vector at the sentence level. W_{sh} and W_{sp} are the weight matrices for hidden state and aspect vector and b_s is bias at the sentence level respectively.

3.6 Cardinal Phrase Guided Hierarchical Attention

We deal with the incongruence detection for the news headlines containing cardinal numbers, therefore the most significant information and cue is the cardinal number itself and neighbouring words. For each headline, we extract a word triplet of form $* : \text{Numerical} - \text{value} : *$, where in place of wildcards, there can be any words., E.g. *Loan 1 million*. We call these word triplets as cardinal phrases. We use these cardinal phrases to drive attention to select salient words and sentences at word level and sentence level correspondingly. To do that, we represent each cardinal phrase CP as the summation of embeddings of all three words of word triplet as:

$$\vec{CP} = f(W_p) + f(W_q) + f(W_r) \quad (7)$$

In a very similar fashion to cardinal POS triplet pattern guided attention, we use \vec{CP} to compute the attention weights at both the word and sentence levels.

$$\alpha_{ij}^c = \frac{\exp(e(HS_{ij}, \vec{CP}))}{\sum_{k=1}^l \exp(e(HS_{ik}, \vec{CP}))} \ \& \ \beta_j^c = \frac{\exp(e(HS_j, \vec{CP}))}{\sum_{k=1}^o \exp(e(HS_k, \vec{CP}))} \quad (8)$$

3.7 Headline Guided Hierarchical Attention

The objective of the headline driven attention is to select words and sentences in the news body text, which are relevant and topically aligned with headline content. The cardinal POS-tag pattern and cardinal phrase carry useful information regarding cardinal values but to capture the whole context of the headline and it's influence on news body text, we can not get rid of headline driven attention. We represent each headline \vec{h} as the summation of embeddings of all the words contained in it as:

$$\vec{h} = \sum_{x=1}^l f(w_x) \quad (9)$$

In a very similar fashion to cardinal POS triplet pattern guided attention, we use \vec{h} to compute the attention weights at both the word and sentence levels.

$$\alpha_{ij}^h = \frac{\exp(e(HS_{ij}, \vec{h}))}{\sum_{k=1}^l \exp(e(HS_{ik}, \vec{h}))} \ \& \ \beta_j^h = \frac{\exp(e(HS_j, \vec{h}))}{\sum_{k=1}^o \exp(e(HS_k, \vec{h}))} \quad (10)$$

3.8 Fusion of Attention Weights and Classification

We compute the overall attention weights from three kinds of attention mechanisms: POS-pattern-driven, Cardinal-phrase-driven, and headline driven attention at both the word and sentence levels. At the word level:

$$\alpha_{i,j} = (\alpha_{i,j}^p + \alpha_{i,j}^c + \alpha_{i,j}^h)/3 \ \& \ S^j = \sum_{i=1}^l \alpha_{ij} HS_{ij} \quad (11)$$

where $\alpha_{i,j}^p$, $\alpha_{i,j}^c$ and $\alpha_{i,j}^h$ are the attention weight vectors from POS-pattern, Cardinal-phrase and headline-attention at the word level. S^j is the formed sentence representation after overall attention for the j^{th} sentence. At the sentence level:

$$\beta_j = (\beta_j^p + \beta_j^c + \beta_j^h)/3 \ \& \ D = \sum_{j=1}^o \beta_j HS_j \quad (12)$$

where β_j^p , β_j^c , and β_j^h are the attention weight vectors from POS-pattern, Cardinal-phrase and headline-attention at the sentence level, and D is the formed document representation after overall attention. The document representation D is used with a Softmax layer with softmax cross-entropy with logits as loss function for the classification. We compute the predicted label \hat{y} as: $\hat{y} = \text{softmax}(W_{cl}D + b_{cl})$. Where W_{cl} and b_{cl} are the weight matrix and bias term.

4 DATASET CREATION

For evaluation, we create the datasets driven by two publicly available datasets, NELA17 and Click-bait Challenge¹ (cf. Table 1). Yoon et al. [30] provide a script² to create the NELA17 dataset from an original news collection NELA17 dataset³. The NELA17 dataset comprises of 45521 congruent and 4551 incongruent news headline-body pairs. The Click-bait Challenge dataset is created via crowd-sourcing based annotation of a collection of social media posts. The Click-bait Challenge dataset contains 16150 and 4883 social media posts, which are annotated as congruent and incongruent correspondingly. Using NELA17 and Click-bait Challenge datasets, we derive two new datasets, in which all the news headlines in NELA17 and all the social media posts in Click-bait Challenge, contain a numerical value. We call these two new datasets as Derived NELA17 dataset and Derived Click-bait Challenge dataset. We create the new datasets using these steps:

- (1) We use POS tagger to get the words in headlines tagged with one of the corresponding Penn Treebank POS Tag Set.
- (2) We keep all the headline-body pairs in which pos-tag CD (cardinal) appears in headline.

The statistics of these datasets are reported in Table 2. We also extract two new features for each headline-body pair:

- A pos-tag triplet of form $* : CD : *$, where in place of stars, there can be JJ, NN, VB etc. E.g. *NN : CD : JJ*. We call this pos-tag triplet as Cardinal POS-tag Pattern. For a vector representation for each of the cardinal pos-tag pattern, we create a trainable embeddings of the size of 100 dimensions and initialize them with uniformly random weights. The weights for these embeddings are learned jointly using hierarchical attention in the POSHAN model.
- A word triplet of form $* : \text{Numerical} - \text{value} : *$, where in place of stars, there can be any words. E.g. *Loan 1 million*. We call this word triplet as Cardinal Phrase.

¹<http://www.clickbait-challenge.org/>

²<https://github.com/sugoiiii/detecting-incongruity-dataset-gen>

³<https://github.com/BenjaminDHorne/NELA2017-Dataset-v1>

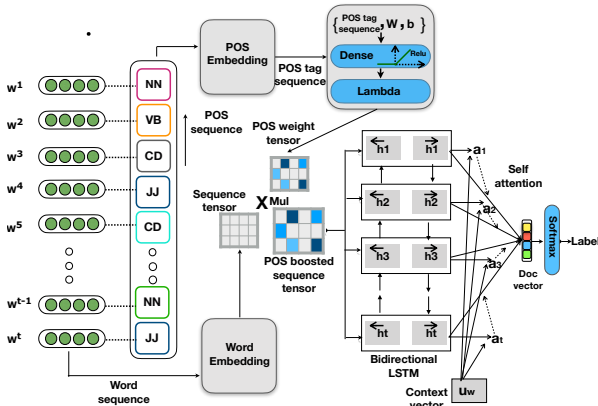


Figure 4: Depiction of the POS tag guided attention model. From left to right, the sequence of words of the news item and corresponding part-of-speech tags go through word embedding and POS embedding layers. Attention scores are computed for each POS tag via a dense layer with ReLU activation. The resultant weight score matrix is multiplied to the word embedding sequence matrix. Attended representation is fed to a Bi-LSTM unit. Lastly, the resultant document representation from Bi-LSTM is used with a softmax layer for classification.

Table 1: Dataset Statistics for NELA17

Statistics	NELA17	Derived NELA17
Incongruent	45521	6234
Congruent	45521	7766
Total	91042	14000

Table 2: Dataset Statistics for Click-bait Challenge

Statistics	Click-bait Challenge	Derived Click-bait Challenge
Incongruent	4883	754
Congruent	16150	2681
Total	21033	3435

5 EXPERIMENTAL EVALUATION

5.1 Experimental Details

5.1.1 *Baselines.* We compare our model with the following baselines:

SVM [7]. : We start with feature-based methods utilizing support vector machine (SVM) by considering both linguistic and statistical features. In specific, we use word tri-grams, four-grams, and part-of-speech bi-grams and tri-grams as features to learn a classifier using SVM. Usually, a click-bait headline contains word phrases like “what happens if” and “You will Never Believe”, which can

be easily captured by tri-grams and four-grams. Besides, POS tag combinations such as “PRP WD RB” are more frequent in incongruent headlines than incongruent ones, therefore, part-of-speech bi-grams and tri-grams are used to learn this distinguishing feature.

LSTM [11]. : We use long short term memory unit to encode both headline and body pair and apply softmax for the classification. We use pretrained GloVe embeddings of size 100 and the size of the hidden states of the Bi-LSTM unit is kept at 200. The concatenation of the news headline and the body text is used as input to the Bi-LSTM based encoder.

POSat. : We propose a baseline method called as POS-tag guided Attention (POSat) and compare the performance of our proposed approach *POSHAN*, as this method uses POS tags to give importance to certain words. This method is inspired by a recent work [27]. We use NLTK POS tagger to tag each word in the news headline/body pairs and maintain the mapping between words and corresponding POS tags using an index. POS tags are categorized into 6 semantic categories for sake of simplicity and brevity.

- (1) **Noun chunk:** NN, NNS, NNP, NNPS
- (2) **Verb chunk:** VB, VBD, VBG, VBN, VBP, VBZ
- (3) **Adjective chunk:** JJ, JJR, JJS
- (4) **Pronoun chunk:** WP, WP
- (5) **Adverb chunk:** WRB
- (6) **Cardinal numbers chunk:** CD

The POS tag for each word is represented as a 6-dimensional vector $g(x_i) \in \mathbb{R}^6$ of length 6. The weights of these embedding vectors are initialized in two ways.

- (1) **Initialize with very less value, close to zero:** In case of all zeros initialization, model performs well.
- (2) **Initialize with random weights:** In case of random weight initialization performance of the model degrades.

These POS tag vector sequences are fed into a fully connected POS tag embeddings layer so that their weights are also trainable. Each part-of-speech category is assigned with one attention weight θ_i , which will be learned during training. In this way, each word is represented by $f'(w_i) = f(w_i) \times \theta_i$. We train a small neural network with only single hidden layer to learn weights for each POS tag category and then we use a custom lambda layer to reshape the POS weight tensor into a compatible shape so that we can boost each word vector with its corresponding POS weight vector.

In very similar fashion to LSTM baseline, pretrained GloVe embeddings of size 100 are used to represent the word vectors and hidden states of the LSTM unit is kept at 200. The nltk library⁴ with MaxEnt POS Tagger [21] is used to tag the concatenation of the news headline and the body text.

Yoon [30]. : This is a state of the art method for news headline incongruence detection. It uses a hierarchical dual encoder based model which uses headline guided attention to learn the contextual representation. The original [30] paper uses a Korean news collection as dataset for evaluation but they also release an English version of the dataset called as NELA17. We use their model with NELA17 dataset, keeping all the settings as prescribed in [30].

⁴<https://pypi.org/project/nltk/>

Table 4: Comparison of the proposed model POSHAN with various state-of-the-art baseline models for click-bait challenge dataset. The results for POSHAN are statistically significant ($p - value = 2.29e^{-3}$ for click-bait challenge dataset using pairwise student’s t-test.

Derived Click-bait challenge Dataset		
Model	Macro F1	AUC.
SVM [7]	0.596	0.608
LSTM [11]	0.604	0.617
POSA	0.614	0.620
BERT-Sent_Pair[8]	0.637	0.649
Yoon [30]	0.646	0.659
MuSeM[17]	0.698	0.717
POSHAN	0.739	0.748
Click-bait challenge Dataset		
Model	Macro F1	AUC.
SVM [7]	0.618	0.629
LSTM [11]	0.630	0.641
POSA	0.636	0.649
BERT-Sent_Pair[8]	0.653	0.662
Yoon [30]	0.660	0.678
MuSeM[17]	0.735	0.747
POSHAN	0.743	0.761

Table 3: Comparison of the proposed model POSHAN with various state-of-the-art baseline models for NELA17 Dataset. The results for POSHAN are statistically significant ($p - value = 1.32e^{-2}$ for NELA17 Dataset using pairwise student’s t-test

Derived NELA17 Dataset		
Model	Macro F1	AUC.
SVM [7]	0.608	0.610
LSTM[11]	0.627	0.639
POSA	0.624	0.637
BERT-Sent_Pair [8]	0.642	0.658
Yoon [30]	0.653	0.659
MuSeM[17]	0.703	0.721
POSHAN	0.748	0.763
Original NELA17 Dataset		
Model	Macro F1	AUC.
SVM [7]	0.622	0.637
LSTM [11]	0.642	0.663
POSA	0.648	0.669
BERT-Sent_Pair[8]	0.677	0.683
Yoon [30]	0.685	0.697
MuSeM[17]	0.752	0.769
POSHAN	0.765	0.783

BERT-Sent_Pair [8]. : We fine tune a pretrained BERT model for sequence pair classification task. We utilize the Hugging face transformers and dataset libraries to download pre-trained model.

We use pre-built "BertForSequenceClassification", provided by Hugging face library. Headlines and body pairs are packed together into a single sequence with adequate padding.

MuSem [17]. : This is a very recent work related to title incongruence detection, which uses both the NELA17 and Click-bait challenge datasets for evaluation. The authors propose a method that uses inter-mutual attention-based semantic matching between the original and a synthetically generated headlines via generative adversarial network based techniques, which utilises the difference between all the pairs of word embeddings of words involved and computes mutual attention score matrix.

5.2 POSHAN Implementation Details

The POSHAN model is implemented using TENSORFLOW 1.10.0⁵ platform. For performance evaluation, Macro F1, and Area Under the ROC Curve (AUC) scores are used as performance metrics. We keep the size of hidden states of bi-directional Long Short-term Units (LSTM) as 300, the size of embedding dimensions of pretrained BERT [8] embeddings as 768. We use softmax cross-entropy with logits as the loss function. We keep the learning rate as 0.003, batch size as 128, and gradient clipping as 6. The parameters are tuned using a grid search. We use 50 epochs for each model and apply early stopping if validation loss does not change for more than 5 epochs. We keep maximum words in a sentence as 45 and maximum number of sentences in a news body text as 35.

⁵<https://www.tensorflow.org/install/source>

Handling the multiple cardinal values. : There are some cases where news headlines contain multiple cardinal values such as in fig 1. At the training time, to utilize the context of all cardinal values present in the headline, we replicate the news headline and body pair for each cardinal value in train set. At the test time however, we concatenate all the learned cardinal POS tag vectors pertaining to the same news headline and use this overall POS tag vector to guide the attention.

Extraction of BERT Embeddings [29]. : We use bert-as-service, which utilizes extract_features.py file from original BERT implementation, to extract the word embeddings from pretrained BERT model. We fine tune uncased_L-12_H-768_A-12 pretrained BERT model for sentence pair classification task. We set pooling_strategy argument to NONE and use our own tokenizer. We use fine tuned BERT model to extract embeddings of 768 dimensions for each word.

5.3 Results

In this section, we compare the results of the POSHAN model with the baselines and state-of-the-art methods.

5.3.1 Results for NELA17 Dataset. In Table 3, we observe that in case of Derived NELA17 dataset, all the deep learning based methods outperform the non-deep learning method such as SVM model, which uses linguistic features and gets 0.608 and 0.610 in terms of Macro F1 and AUC. The POSA model with Macro F1 score as 0.624 and AUC as 0.637, performs comparable with vanilla LSTM model. In our experiments, we introspect that the design decision in POSA model to apply POS-tag guided attention at the POS-tag

Table 5: Ablation study of POSHAN and Yoon[30] model conducted on derived NELA 17 Dataset.

Derived NELA 17 Dataset		
Scenario	Macro F1	AUC.
Original POSHAN	0.748	0.763
1) Remove Cardinal POS Att	0.726	0.742
2) Remove Cardinal Phrase Att	0.731	0.749
3) Replace Headline Att with Headline Enc	0.648	0.669
4) Replace BERT with Glove	0.716	0.736
5) Replace Bi-LSTM with Bi-GRU	0.746	0.761
6) Replace Bi-LSTM with LSTM	0.741	0.759

Derived NELA 17 Dataset		
Scenario	Macro F1	AUC.
Original Yoon	0.653	0.659
1) Remove Headline Att	0.593	0.595
2) Replace para to sent level Att	0.649	0.653
3) Replace Glove with W2V	0.610	0.618
4) Replace Bi-LSTM with Bi-GRU	0.652	0.657
5) Replace Bi-LSTM with LSTM	0.641	0.648

chunk level, does not result in effective representation and provides very less intended effects of POS types on words. This way of POS-tag guided attention learns the attention score at the POS category level only as discussed in Section 5.1.1 such as Noun chunk, Verb chunk etc.

The BERT-Sent_Pair with Macro F1 as 0.642 and AUC as 0.658, outperforms the POSAt model with significant difference. This gain can be attributed to the better contextual representation of words, learned in form of transformer based BERT embeddings. On the other hand, Yoon model [30] performs slightly better than BERT-Sent_Pair with Macro F1 as 0.653 and AUC as 0.659. In addition to hierarchical encoder, which captures the complex structure of the news body content, having inherent hierarchical nature, Yoon model also uses a headline driven hierarchical attention, which not only selects salient and relevant words and sentences but also reduces the effective length of the news body. In contrast, vanilla LSTM, POSAt and even BERT-Sent_Pair model did not scale well for long text sequences. The MuSem model uses generative adversarial network based synthetic headline generation methods to generate a very low dimensional headline corresponding to news body and applies a novel mutual attention based semantic matching for incongruence detection. The MuSem model achieves significant gains over Yoon model, due to low dimensional representation of news body and effective semantic matching technique. The proposed POSHAN model beats all the other methods achieving 0.748 and 0.763 as Macro F1 and AUC, respectively. The potential reason behind this better performance is superior document representation learned due to proposed attention mechanisms, which give adequate importance to significant cardinal values present in headline. In contrast, both Yoon and MuSem models fail to capture cues pertaining to cardinal patterns and phrases.

In case of Original NELA17 Dataset, we notice a very similar trend as with derived NELA17 dataset, however performance of

all the models improved with a significant margin. On the other hand, POSHAN happens to yield more improvement in performance compared to other models for original dataset as a bonus. These gains can be attributed to cardinal POS-tag pattern based attention and cardinal phrase guided attention in addition to headline guided attention significantly.

5.3.2 Results for Click-bait Challenge Dataset. In case of both the derived and original click-bait challenge dataset also, we see a very similar performance chart. All the deep learning based methods outperform the non-deep learning method such as SVM model With 0.596 and 0.608 in Macro F1 and AUC. The MuSem model with 0.698 and 0.717 in Macro F1 and AUC, outperforms all the other baselines. The proposed model POSHAN performs better than MuSem model with significant gains and these gains can be explained by very similar reasoning, as provided in Section 5.3.1.

5.4 Ablation Study

In Table 5, we report an ablation study of POSHAN using Derived NELA 17 Dataset. We used derived dataset for ablation study rather than original dataset because we want to assess the importance of different components of the POSHAN with major focus of this paper, which is news headlines with important numerical values. In the ablation version 1), we remove the cardinal POS-tag pattern guided attention and keep the other two methods of attention intact and this step results in significant decrease in performance, which proves the usefulness and effectiveness of the cardinal POS-tag pattern guided attention. This corroborates with our original hypothesis and intuition. In the ablation version 2), we remove cardinal phrase attention and observe very similar decrease in performance. In the ablation version 3), we replace headline guided attention with headline encoder, in which we encode the words of news headlines in addition to news body words and concatenate the overall encoded sequence. We observe that without headline guided attention, model performs poorly because just a simple concatenation of encoded body and headline word sequences does not result in contextually important representation. We can conclude from 1), 2) and 3), that although all the three attention mechanism are effective individually too but combination of all the three becomes more effective. In the ablation version 4), we replace the pretrained BERT embeddings with GloVe [15], due to which the performance degrades drastically. The reason behind such a drop in the results is that the BERT embeddings provide superior contextual information than GloVe pretrained embeddings. We do not see much change in results in ablation version 5) as Bi-GRU [6] and Bi-LSTM perform pretty much the same with our dataset. The performance of the model decreases a bit with replacement of Bi-LSTM with LSTM units in ablation version 6) and the obvious reason behind this better context learned by Bi-LSTM compared to LSTM units.

5.5 Error Analysis

We conduct an error analysis of MuSem [17] and POSHAN model with Derived Click-bait challenge dataset in Table 6. In the case of MuSem model, we observe 235 false negatives(FN) and 201 false positives (FP), on the other hand, POSHAN produces 207 false negatives and 179 false positives. We notice that the major improvement

Headline:- Immigration Expert : US Will Have 100 Million New Immigrants in Next 50 Years.

Yoon Model

These projections show that new immigrants and their descendants will drive most U.S. population growth in the coming 50 years, as they have for the past half-century. Among the projected 441 million Americans in 2065, 78 million will be: immigrants and 81 million will be people born in U.S. to immigrant parents.

POSHAN Model

These projections show that new immigrants and their descendants will drive most U.S. population growth in the coming 50 years, as they have for the past half-century. Among the projected 441 million Americans in 2065, 78 million will be: immigrants and 81 million will be people born in U.S. to immigrant parents.

Figure 6: Attention weight visualization: Word level attention weights from Yoon Model and POSHAN model for an anecdotal example are presented by highlighting the individual words (Best viewed in color). The depth of the color represents the strength of the attention weights.

Table 6: Error Analysis of Yoon model [30] and POSHAN with Derived Click-bait challenge Dataset

Model	False Positives	False Negatives
MuSem	201	235
POSHAN	179	207

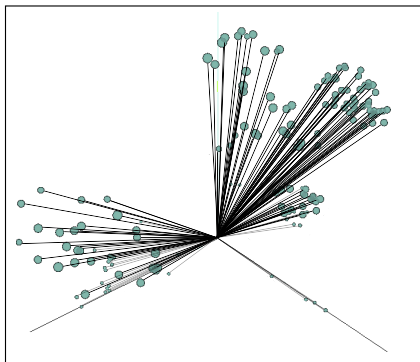


Figure 5: Visualization of Cardinal POS Pattern Embeddings with POSHAN model, occurs in false negatives from 235 to 207, and most of these incorrectly predicted samples were related to important cardinal figures mentioned in the news headlines such as ‘Indiana couple admits to stealing 1.2 Million dollars from Amazon’. We also observe some incorrectly predicted false-positive cases by POSHAN model because of the wrong POS-tag assignment by POS tagger and due to this POSHAN misses out on the opportunity to consider those cardinal POS patterns and cardinal phrases.

5.6 Visualization of Cardinal POS Pattern Embeddings

In the Figure 5, we present a visualization of cardinal pos-tag patterns. To visualize the learned embeddings of the cardinal pos-tag

patterns, we use t-Distributed Stochastic Neighbor Embedding (t-SNE) [25] with parameters as perplexity = 10, learning rate = 0.1 and iterations = 1000. The t-SNE method produces the visualization in a low dimensional space. We observe that the Cardinal POS Patterns have formed clearly separated clusters in embedding space, which connotes the congruence and incongruence labels. We also observe that cardinal POS patterns with similar tags such as (NN : CD : JJ) and (NNS : CD : JJ) are closer in embedding space and on the other hand the patterns with disjoint tag combinations such as (NN : CD : JJ) and (VBG : CD : CD) are farther apart from each other.

5.7 Visualization of Attention Weights

In Figure 6, to analyse the interpretability of our model POSHAN and to showcase the effectiveness of the proposed attention mechanism in forming the contextually important representations, we visualize the attention maps and compare it with Yoon model. In Figure 6, we use distribution of word level attention weights learned from both POSHAN and Yoon model for an anecdotal example by highlighting the individual words. The depth of the color highlights represents the distribution of attention weights. Despite of common headline driven attention in both the models, we observe some clear differences between attention maps of Yoon model and POSHAN model due to additional cardinal pos-tag pattern and cardinal phrase guided attention mechanisms in POSHAN model. The Yoon model successfully attends some words such as ‘immigrants’, ‘growth’ and ‘projections’ etc. relevant to headline context but fails to capture any words pertaining to significant cardinal phrases such as ‘78 million’ and ‘50 years’ etc. On the other hand, POSHAN model not only gives the importance to the words captured by Yoon model but also, it focuses on important cardinal phrases, which is in concert with our intuition about modeling the POS-tag pattern and cardinal phrase based attention.

6 CONCLUSIONS AND FUTURE WORK

In this paper, we introduce a novel task of incongruence detection in the news when the news headline contains cardinal values. The existing methods fare poorly as they fail to capture the context, pertaining to cardinal values. We present a joint neural model POSHAN, which uses the fine-tuned BERT embeddings with three kinds of hierarchical attention mechanisms, namely cardinal POS-tag pattern guided, cardinal phrase guided and news headline guided attention. In the ablation study, we found that cardinal POS-tag pattern guided attention is very significant and effective in forming the cardinal quantity informed document representation. In the evaluation with two publicly available datasets, we notice that POSHAN outperforms all the baselines and state-of-the-art methods. Visualization of cardinal POS-tag pattern embeddings and overall attention weights establish the effectiveness of the proposed model, decipher the model’s decisions and make it more interpretable and transparent. In the future, we plan to model the degree of importance of cardinal values in news headlines and also we envisage an assessment of the applicability of the proposed model in case of textual entailment and fact verification tasks such as FEVER [24] dataset, in presence of cardinal values.

REFERENCES

- [1] Ankesh Anand, Tanmoy Chakraborty, and Noseong Park. 2017. We Used Neural Networks to Detect Clickbait: You Won't Believe What Happened Next!. In *Advances in Information Retrieval*. Springer International Publishing, 541–547.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. [n. d.]. Neural Machine Translation by Jointly Learning to Align and Translate. In *ICLR 2015 as oral presentation*.
- [3] Prakhar Biyani, Kostas Tsioutsoulis, and John Blackmer. 2016. “8 Amazing Secrets for Getting More Clicks”: Detecting Clickbaits in News Streams Using Article Informality. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence (AAAI'16)*. AAAI Press, 94–100.
- [4] A. Chakraborty, B. Paranjape, S. Kakarla, and N. Ganguly. 2016. Stop Clickbait: Detecting and preventing clickbaits in online news media. In *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*. 9–16.
- [5] Yimin Chen, Niall J. Conroy, and Victoria L. Rubin. 2015. Misleading Online Content: Recognizing Clickbait as “False News”. In *Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection (WMDD '15)*. Association for Computing Machinery, New York, NY, USA, 15–19. <https://doi.org/10.1145/2823465.2823467>
- [6] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1724–1734. <https://doi.org/10.3115/v1/D14-1179>
- [7] Corinna Cortes and Vladimir Vapnik. 1995. Support-Vector Networks. *Mach. Learn.* 20, 3 (Sept. 1995), 273–297. <https://doi.org/10.1023/A:1022627411411>
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [9] 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*. Association for Computational Linguistics, San Diego, California, 1163–1168. <https://doi.org/10.18653/v1/N16-1138>
- [10] Maksym Gabielkov, Arthi Ramachandran, Augustin Chaintreau, and Arnaud Legout. 2016. Social Clicks: What and Who Gets Read on Twitter?. In *Proceedings of the 2016 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Science (SIGMETRICS '16)*. Association for Computing Machinery, New York, NY, USA, 179–192. <https://doi.org/10.1145/2896377.2901462>
- [11] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Comput.* 9, 8 (Nov. 1997), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [12] John S. Justeson and Slava M. Katz. 1995. Technical terminology: some linguistic properties and an algorithm for identification in text. *Natural Language Engineering* 1, 1 (1995), 9–27. <https://doi.org/10.1017/S1351324900000048>
- [13] Vaibhav Kumar, Dhruv Khatwar, Siddhartha Gairola, Yash Kumar Lal, and Vasudeva Varma. 2018. Identifying Clickbait: A Multi-Strategy Approach Using Neural Networks. In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '18)*. Association for Computing Machinery, New York, NY, USA, 1225–1228. <https://doi.org/10.1145/3209978.3210144>
- [14] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. [arXiv:cs.CL/1907.11692](https://arxiv.org/abs/1907.11692)
- [15] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Eds.). Curran Associates, Inc., 3111–3119. <http://papers.nips.cc/>
- [16] Rahul Mishra and Vinay Setty. 2019. SADHAN: Hierarchical Attention Networks to Learn Latent Aspect Embeddings for Fake News Detection (*ICTIR '19*). 197–204.
- [17] Rahul Mishra, Piyush Yadav, Remi Calizzano, and Markus Leippold. 2020. MuSeM: Detecting Incongruent News Headlines using Mutual Attentive Semantic Matching. In *International Conference on Machine Learning and Applications (ICMLA) 2020*. Miami, Florida.
- [18] Paul Neculou, Maarten Versteegh, and Mihai Rotaru. 2016. Learning Text Similarity with Siamese Recurrent Networks. In *Proceedings of the 1st Workshop on Representation Learning for NLP*. Association for Computational Linguistics, Berlin, Germany, 148–157. <https://doi.org/10.18653/v1/W16-1617>
- [19] Rahul Potharaju, Navendu Jain, and Cristina Nita-Rotaru. 2013. Juggling the Jigsaw: Towards Automated Problem Inference from Network Trouble Tickets. In *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*. USENIX Association, Lombard, IL, 127–141. <https://www.usenix.org/conference/nsdi13/technical-sessions/presentation/potharaju>
- [20] Martin Potthast, Sebastian Köppl, Benno Stein, and Matthias Hagen. 2016. Clickbait Detection. *ECIR '16* 1 (2016).
- [21] Adwait Ratnaparkhi. 1996. A Maximum Entropy Model for Part-Of-Speech Tagging. In *Conference on Empirical Methods in Natural Language Processing*. <https://www.aclweb.org/anthology/W96-0213>
- [22] Julio Reis, Pedro Olmo, Raquel Prates, Haewoon Kwak, and Jisun An. [n. d.]. Breaking the News : First Impressions Matter on Online News. ([n. d.]), 357–366.
- [23] K. Shu, S. Wang, T. Le, D. Lee, and H. Liu. 2018. Deep Headline Generation for Clickbait Detection. In *2018 IEEE International Conference on Data Mining (ICDM)*. 467–476.
- [24] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. Association for Computational Linguistics, New Orleans, Louisiana, 809–819. <https://doi.org/10.18653/v1/N18-1074>
- [25] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE.
- [26] Zhiguo Wang, Wael Hamza, and Radu Florian. 2017. Bilateral Multi-Perspective Matching for Natural Language Sentences. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*. 4144–4150. <https://doi.org/10.24963/ijcai.2017/579>
- [27] Z. Wang, X. Liu, L. Wang, Y. Qiao, X. Xie, and C. Fowlkes. 2018. Structured Triplet Learning with POS-Tag Guided Attention for Visual Question Answering. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. 1888–1896.
- [28] Wei Wei and Xiaojun Wan. 2017. Learning to Identify Ambiguous and Misleading News Headlines. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI'17)*. AAAI Press, 4172–4178.
- [29] Han Xiao. 2018. bert-as-service. <https://github.com/hanxiao/bert-as-service>.
- [30] Seunghyun Yoon, Kunwoo Park, Joongbo Shin, Hongjun Lim, Seungpil Won, Meeyoung Cha, and Kyomin Jung. [n. d.]. Detecting Incongruity Between News Headline and Body Text via a Deep Hierarchical Encoder. ([n. d.]). [arXiv:arXiv:1811.07066v2](https://arxiv.org/abs/1811.07066v2)