

# An Attention-Based Neural Network Using Human Semantic Knowledge and Its Application to Clickbait Detection

FENG WEI  (Student Member, IEEE), AND UYEN TRANG NGUYEN (Member, IEEE)

Department of Electrical Engineering and Computer Science, York University, Toronto, Ontario M3J 1P3, Canada

CORRESPONDING AUTHOR: FENG WEI (e-mail: fwei@eecs.yorku.ca)

This work was supported in part by a Discovery Grant funded by the Natural Sciences and Engineering Research Council of Canada.

**ABSTRACT** Clickbait is a commonly used social engineering technique to carry out phishing attacks, illegitimate marketing, and dissemination of disinformation. As a result, clickbait detection has become a popular research topic in recent years due to the prevalence of clickbait on the web and social media. In this article, we propose a novel attention-based neural network for the task of clickbait detection. To the best of our knowledge, our work is the first that incorporates human semantic knowledge into an artificial neural network, and uses linguistic knowledge graphs to guide attention mechanisms for the clickbait detection task. Extensive experimental results show that the proposed model outperforms existing state-of-the-art clickbait classifiers, even when training data is limited. The proposed model also performs better or comparably to powerful pretrained models, namely, BERT, RoBERTa, and XLNet, while being much more lightweight. Furthermore, we conducted experiments to demonstrate that the use of human semantic knowledge can significantly enhance the performance of pretrained models in the semisupervised domain such as BERT, RoBERTa, and XLNet.

**INDEX TERMS** Clickbait detection, fake news, human semantic knowledge, knowledge base, neural networks.

## I. INTRODUCTION

Clickbait is designed to entice Internet users into clicking on a hyperlink that leads to a web page or another piece of online content such as an article, an image or a video. There are three principal uses of clickbait:

- 1) *Phishing attacks*: Attackers send malicious messages to victims that trick them into clicking on a hyperlink. The purpose is to steal sensitive information such as login credentials, banking information or credit card numbers, or to install malicious software such as ransomware on the victims' computers.
- 2) *Legitimate and illegitimate marketing*: Marketers use attention-grabbing headlines to entice users into clicking on hyperlinks that lead to online articles or advertisements. Marketers make money based on the number of clicks that an article or advertisement receives.
- 3) *Spreading misinformation and disinformation (fake news)*: Similarly, attention-grabbing headlines are used

to build suspense and sensation in order to lure users into clicking on hyperlinks that lead to articles or online content with fake news.

The focus of our work in this article is on cases (2) and (3). In particular, we determine whether a headline is clickbait or not by matching its content to the textual content of the article linked to it. If the content of the headline closely matches that of the linked article, it is not clickbait. Otherwise it is classified as clickbait. An example is given in Table 1.

It should be noted that not all articles with clickbait headlines are fake news. However, clickbait headlines are one of the hallmarks of fake news. Therefore, the classification of clickbait is useful as the first step toward a more fine-grained classification of fake news (e.g., by examining the content of the article for authenticity).

Clickbait detection has become a popular research topic in the past few years due to the prevalence of clickbait on the Internet. Chen et al. [1] examined both contextual (e.g., semantic

**TABLE 1.** Examples of Clickbait and Non-Clickbait From the FNC Challenge dataset [3]

◆ <b>Headline</b>	EXCLUSIVE: Apple To Unveil The Long-Awaited Retina MacBook Air At Its “Spring Forward” Event
♣ <b>Body</b>	Last week, Apple sent out the invites for its “Spring Forward” event, . . . The Retina MacBook Air, currently codenamed “MacBook Stealth” internally, will also come with a modest price decrease according to previous reports.
♠ <b>Candidate</b>	[clickbait, non-clickbait]
◆ <b>Headline</b>	Next-generation Apple iPhones’ features leaked
♣ <b>Body</b>	When faced with the choice of feasting on a fine meal of human while listening to Justin Bieber’s music . . . his cellphone rang, and the ringtone of Justin Bieber hit “Baby” startled the animal . . . “I know that sort of ringtone isn’t to everyone’s taste but my granddaughter loaded it onto my phone for a joke.” . . .
♠ <b>Candidate</b>	[clickbait, non-clickbait]

and lexical features) and non-contextual information (e.g., figures, user behavior). The authors applied conventional machine learning methods such as support vector machines (SVM) and naïve Bayes to these features for final predictions. Potthast et al. [2] examined the features from both the headlines and linked articles/webpages. Additionally, the authors leveraged linguistic information (e.g., the average word length and sentiment preference) and lateral information (e.g., the authorship of a headline). They then fed a combination of these features into conventional machine learning models such as random forest, naïve Bayes, and logistic regression.

Recently, artificial neural networks (ANN) have been widely used to detect clickbait [4], [5], due to their effectiveness and efficiency in deriving complex associations from a dataset [6] and processing large volumes of features. For example, most top-ranked competitors in the 2017 Clickbait Challenge [7] used deep learning in their clickbait detection models. The first ranked team [5] applied bi-directional gated recurrent units (BiGRU) with incorporated attention mechanisms [8] to model contextual information from the clickbait dataset. In addition, Glenski et al. [4] adopted long short term memory (LSTM) [9] and convolutional neural networks (CNN) to capture both textual and visual information from the dataset and feed them into the models.

However, the above ANN-based models [4], [5] mainly focus on integrating a variety of neural network structures. The importance of semantic correlations between a headline and its linked article was overlooked in the above works. In particular, every pair of words, one from the headline and the other from the linked article (headline-article word pair), can be examined to derive semantic correlations between them and, ultimately, between the headline and the article [10]. The semantic correlations between word pairs are then used to build linguistic knowledge graphs [11]; the graphs will enable

An example from FNC Challenge dataset.

**Headline:** Justin Bieber *ringtone* scares bear, saves Russian *fisherman*  
**Body:** A Russian *man* was recently able to fend off a bear attack with a Justin Bieber *song*.  
**Label:** *non-clickbait*

**FIGURE 1.** An example illustrating the importance of human semantic knowledge to clickbait detection. Without human semantic knowledge, this example will be incorrectly classified as clickbait. In the example, the proposed model can find the correct answer because it knows “ringtone” is a synonym of “song,” and “fisherman” is a hyponym of “man” thanks to the use of human semantic knowledge.

neural models to identify similarity patterns in a latent feature space more effectively.

Only a few existing models have used the correlations between headlines and their linked articles for clickbait identification. Biyani et al. [12] used similarities between the headline and the first five sentences of the linked article for clickbait detection, integrated with conventional handcrafted textual features (e.g., the number of upper case words and the number of acronyms). Kumar et al. [13] used Siamese networks to assess similarities between the textual content and visual images. The authors fed the similarity features into several fully connected neural network layers. Dong et al. [14] used bi-directional gated recurrent units (GRU) to learn contextual features of headlines and their linked articles. Additionally, the authors applied a cosine similarity-based method to learn global matching similarities between headlines and their linked articles.

There are, however, two shortcomings in the above clickbait identification models. First, they do not leverage human semantic knowledge, which is essential to determine semantic correlations between a headline and its linked article for the clickbait detection task [10]. The example in Fig. 1 illustrates the importance of human semantic knowledge to clickbait detection. Second, the methods they use to match a headline with its linked article, such as cosine similarity [12] or simple attention-based networks [14], do not adequately capture similarity information.

To overcome these challenges, we propose an attention-based neural network using human semantic knowledge, which can be applied to clickbait detection and several other tasks (e.g., word sense disambiguation and machine reading comprehension). We introduce two novel mechanisms in the proposed network. First, we incorporate human semantic knowledge into the proposed neural network, which is inspired by the work in neuroscience by Chen et al. [15]. We use WordNet [16], a lexical database of semantic relations between words, to acquire semantic correlations of headline-article word pairs. Intuitively, semantic knowledge of word pairs such as hyponyms, hypernyms, and/or synonyms can facilitate the similarity matching of words in a headline to words in the linked article. Therefore, for each headline and its linked article, the proposed model generates a linguistic knowledge graph using WordNet that represents semantic correlations of headline-article word pairs. Secondly, the

attention mechanisms, joint attention and self-attention, used by the neural network are enhanced with human semantic knowledge. Specifically, the attention mechanisms use the knowledge graph generated earlier and focus on the most important parts of the graph in order to extract most meaningful headline-article word pairs. For this reason, we say that the proposed attention mechanisms are *knowledge enhanced*.

To the best of our knowledge, our work is the first that

- incorporates human semantic knowledge into a neural network for clickbait detection;
- uses linguistic knowledge graphs built from WordNet to guide the attention mechanisms for the clickbait detection task.

Our contributions are as follows.

- We propose a novel attention-based neural network model named Knowledge-Enhanced Clickbait Detector (KED) that uses linguistic knowledge graphs built from WordNet to guide the attention mechanisms. The proposed neural network can effectively capture discriminative features from local and global similarities. Global similarities are captured using the proposed knowledge-enhanced joint-attention mechanism. To minimize the impact of noise, the model selects the most useful similarity features for the final predictions using the proposed knowledge-enhanced self-attention mechanism.
- We incorporate human semantic knowledge into the neural network and its attention mechanisms to better capture semantic correlations of headline-article word pairs. Experimental results show that this novelty significantly enhances the performance of KED and that of several pretrained models such as BERT, RoBERTa, and XLNet.
- We conducted extensive experiments to evaluate the proposed model on two real-world clickbait datasets, Clickbait Challenge [7] and FNC Challenge [3]. Experimental results show that our model significantly outperforms existing state-of-the-art models/systems [4], [5], [13], [14], [17], [18], [19], [20], [21], even when training data is limited. The proposed model also performs better or comparably to powerful pretrained models, namely, BERT, RoBERTa, and XLNet, while being much more lightweight.
- Based on extensive comparative experiments, we carry out comprehensive analyses which will benefit future studies.

Illegitimate marketing has been used to defraud consumers via products and services such as counterfeit goods, and pyramid, Ponzi, “make money quick” and pump-and-dump schemes. Fake news has had profound impacts on a country’s political and social stability [22], [23], democracy [24], financial markets [25], [26], and public health [27]. Our work in this article, clickbait detection, is a first step towards developing effective countermeasures to the threats of consumer frauds and fake news.

The remainder of this article is organized as follows. Section II presents related work. Section III describes in

detail the proposed attention-based neural network that uses human semantic knowledge. Section IV discusses the application of the proposed neural network to clickbait detection. In Section V, we present and analyze experimental results to validate the effectiveness of pretrained word embeddings, knowledge enhancement to the attention mechanisms, and human semantic knowledge in the proposed model KED. We also compare the performance of the proposed model with that of existing state-of-the-art models/systems. Section VI presents concluding remarks and future work.

## II. RELATED WORK

We briefly review neural networks in Sections II-A, and discuss existing work on clickbait detection in Section II-B.

### A. NEURAL NETWORKS

Artificial neural networks (ANN), with an excellent capacity for automated recognition of patterns and regularities in data, are meant to emulate human brains. They are composed of densely interrelated neurons. After being fed with data, the neurons enable a model to learn the hidden features and produce predictions [28].

A theoretical paradigm called threshold logic [29] can be seen as an initial archetype of ANN. Their study revealed that ANN could be a type of possible instrument for the implementation of artificial intelligence.

Later, an idea of the perceptron was introduced [30]. It consists of an input layer and an output layer, which are linearly connected. Afterward, the popular backpropagation method [31] enables the training of numerous layers ANN like multi-layer perceptron (MLP). In addition, MLPs have been shown to be able to simulate most of the potential functions in practice by using non-linear activation functions (e.g., ReLU) [32].

Due to the shortage of annotated data and insufficient computing performance, however, it was not possible to have large-scale training on ANN (e.g., RNN and CNN). Recently, the advent of fast graphics processing units (GPUs) resolved this issue, which results in 10 or 20 times speed-up for training neural networks [6]. Furthermore, a number of large data sets such as BooksCorpus [33] and ImageNet [34] are published, which successfully benefit the training of deep neural networks. As a result, in a wide range of fields, such as computer vision and natural language processing, ANNs continue to improve the classification performance thanks to the increases in network depths and widths [28].

In recent years, deep learning has contributed to the advancements of numerous fields such as cybersecurity, natural language processing, and computer vision. Since very little manual engineering is required, deep neural networks are able to benefit from the increase in data volume and computing power [6]. In other words, patterns of latent features from a large data set can be effectively and efficiently derived by deep learning techniques. Therefore, more advancements in many applications can be expected in the future thanks to deep learning [28].

## B. EXISTING WORK ON CLICKBAIT DETECTION

As a new topic of study, prior research works on clickbait detection extracted and applied hidden features from the data [1], [2], [35]. In such works, both non-textual information (e.g., images) and textual content were taken into account. In addition to the features of semantic and lexical levels extracted from the content, user behavior characteristics were also considered. A variety of machine learning methods were applied based on those obtained features, such as SVM and naïve Bayes.

Recent research on clickbait detection leverages supervised learning frameworks using a wide range of features, such as readability, forward references, term frequencies [12]), which are linguistic-based; webpage links [2], which is non-linguistic-based; headline stance [36] and user interests [37], [38]. Other studies showed that clickbait could be identified by readability, cardinal numbers, particular adjectives, and nouns describing sentiment, which reflect information on authority and sensationalism [39]. More recently, in order to avoid the labour-intensive task of feature engineering, deep learning has been used more widely for clickbait detection [5], [40], [41].

Razaque et al. proposed a RNN model to determine if a link (URL) in a clickbait message is malicious or harmless [42]. Another RNN model by the same authors determines if the content pointed to by a link (e.g., an article or a social media post) is malicious or harmless [43]. Zhou et al. proposed a clickbait detection model based on graph convolutional networks, and evaluated the model using a dataset in the Chinese language [44]. Recent deep learning models for clickbait detection [20], [21], [45] adapted and fine-tuned transformer models such as BERT [46], RoBERTa [47], XLNet [48], ELECTRA [49] and ALBERT [50]. Mareddy et al. [45] evaluated their transformer-based model using a large dataset in the Telugu language.

To benefit from deep learning techniques, word vectors [51], [52] have also been used to represent textual information in many studies [53]. In [19], the headlines were mapped into word embeddings matrices, which were then fed into convolutional neural networks (CNN). Additionally, recurrent neural networks (RNN), because of their effectiveness in handling sequential data, are commonly adopted for the detection of clickbait. In the Clickbait Challenge competition, all the top five teams applied RNN-based models such as attention-based bi-GRU [5] and LSTM [4] to the analysis of textual input.

Only a few works have studied the correlations between headlines and their linked articles for clickbait identification. Biyani et al. [12] adopted n-gram-based gradient boost decision trees (GBDT) to evaluate information similarities, as well as various measurements such as Coleman-Liau scores [54] and RIX and LIX indices [55]. Kumar et al. [13] applied Siamese networks to determine the textual similarity between headlines and linked articles. The author also considered the similarity between textual descriptions and images. Finally, both textual and visual similarity features were combined to be input into the model for training and inference.

The proposed neural model KED is different from existing clickbait detection models in that it uses

- human semantic knowledge pre-extracted from WordNet to assist in the matching of each given headline and its linked article;
- knowledge graphs to guide the attention mechanisms to extract the most meaningful headline-article word pairs.

As a result, the proposed model performs significantly better than existing state-of-the-art models [5], [14].

An earlier version of our neural network model and preliminary results were first reported in [10]. The earlier model (which we name  $KED_R$  in this article) and results in [10] were based on randomly initialized word embeddings. In this article, we extend the model and optimize its performance by incorporating pretrained word embeddings [56] into KED. In addition, we conducted new experiments and provided new results that compare the performance of KED and with that of current state-of-the-art models when training data is scarce (the case of data scarcity). Another major extension is new experimental results that show that the use of human semantic knowledge significantly improves the classification performance of pretrained models in the semisupervised domain such as BERT, RoBERTa, and XLNet.

## III. THE PROPOSED ATTENTION-BASED NEURAL NETWORK

We propose an attention-based neural network named Knowledge-Enhanced Clickbait Detector (KED) that comprises two parts: (a) human semantic knowledge extraction, and (b) attention mechanisms guided by knowledge graphs.

### A. HUMAN SEMANTIC KNOWLEDGE EXTRACTION USING WORDNET

In this section, we describe in detail our human semantic knowledge extraction scheme, which uses WordNet [16] to extract human semantic knowledge from headline-article word pairs in the dataset. Semantic relations extracted from WordNet are used to build linguistic knowledge graphs, which are then input into the model.

Intuitively, knowledge about antonyms, co-hyponyms, synonyms, hyponyms, hypernyms, and synonyms between given words can potentially help to align word pairs between a given headline and its linked article. For instance, knowledge about synonyms, hypernyms and hyponyms is helpful to identify entailment relations; knowledge about antonyms and co-hyponyms helps to recognize contradiction relations [11].

We extend [11] to extract relations of lexical pairs in a configurable manner by using a hyperparameter  $\tau \in \mathbb{N}$  as follows.

- Synonym relation: If the words in the pair are synonyms in WordNet, the pair yields 1 (true). Otherwise, it yields 0 (false).
- Antonym relation: If the words in the pair are antonyms in WordNet, the pair yields 1 (true). Otherwise, it yields 0 (false).



- **Hypernym relation:** If one word is a hypernym of another word in WordNet, the pair yields value  $1 - n/\tau$ , where  $n$  denotes the number of edges between the two words in a hierarchy, and the hyper-parameter  $\tau$  denotes the exact path length. Otherwise, it yields 0 (false). Consider two word pairs: “fisherman, skilled worker” and “fisherman, worker,” and  $\tau = 5$ . Since ‘skilled worker’ is the direct hypernym of ‘fisherman,’ and ‘worker’ is the direct hypernym of ‘skilled worker’ (i.e., worker  $\xleftarrow[\text{hypernym}]{\text{direct}}$  skilled worker  $\xleftarrow[\text{hypernym}]{\text{direct}}$  fisherman), they yield  $\{\text{fisherman, skilled worker}\} = 1 - \frac{1}{5} = 0.8$  and  $\{\text{fisherman, worker}\} = 1 - \frac{2}{5} = 0.6$ , respectively. It is worth noting that a large  $\tau$  increases the chances of trivial hypernym relations retrieved from the knowledge base, eventually leading to lower performance. On the other hand, a small  $\tau$  would limit the number of hypernym relations obtained from the knowledge base, also resulting in lower performance.
- (d) **Hyponym relation:** It is the inverse of the hypernym feature.
- (e) **Co-hyponym relation:** If the two words have the same hypernym, but they do not belong to the same synset, the pair yields 1 (true). Otherwise, it yields 0 (false).

Note that hyperparameter  $\tau$  enables effective extractions of hypernym relations (and, consequently, hyponym and co-hyponym relations) from the knowledge base by avoiding extracting too few or too many hypernym relations.

The above extracted relations from the knowledge base WordNet applied to the joint-attention mechanism facilitates the matching between headlines and their linked articles.

## B. THE PROPOSED KNOWLEDGE ENHANCED ATTENTION MECHANISMS

Attention mechanisms are essential components of our model. In this section, we describe in detail the proposed knowledge enhanced attention mechanisms.

Attention mechanisms aim to fuse the associated representations of headline-article word pairs. There are two attention mechanisms: joint attention [57] and self attention [58]. Given a headline and its linked article, the purpose of *joint attention* is to fuse the headline representation into the article representation in order to obtain a headline-aware representation of the article [14]. The purpose of *self attention* is to fuse the headline-aware representation of the article into itself in order to obtain the final article representation [5].

Existing works [5], [14] apply either standard dot-scale attention or cosine similarity-aided attention to align a headline and its linked article, without using any external knowledge. In contrast, our proposed neural network leverages WordNet to build knowledge graphs which are then used to guide the focus of the two attention mechanisms. We name the attention mechanisms in the proposed model *knowledge-enhanced joint-attention* and *knowledge-enhanced self-attention*, respectively.

### 1) KNOWLEDGE-ENHANCED JOINT-ATTENTION

Joint attention aims to interconnect the representations of a headline and its linked article, and produces a set of headline-aware feature vectors for the words in the linked article [57]. The input to the knowledge-enhanced joint-attention layer includes the representations of a headline and its linked article from the previous layer of a model, and the output is headline-aware vector representations of the words in the linked articles.

Formally, given a set of linked articles  $B \in \mathbb{R}^n$  with  $n$  words and a set of headlines  $H \in \mathbb{R}^m$  with  $m$  words, let  $E_B \in \mathbb{R}^{d \times n}$  and  $E_H \in \mathbb{R}^{d \times m} \in \mathbb{R}^{n \times m}$  denote the embeddings of the linked articles and headlines, respectively. Following [59], we compute the similarity of each headline-article word pair, rendering a similarity matrix  $J \in \mathbb{R}^{n \times m}$ . Then, each row of  $J$  is normalized by applying the softmax function, and a matrix  $\bar{J}$  is obtained. The article-to-headline attention is then computed as  $\bar{J} \cdot H^T \in \mathbb{R}^{d \times n}$ . The headline-aware vector  $p_i \in P$  of article  $b_i \in B$  is computed for the words in  $b_i$  as follows [57], [60]:

$$p_i = \mathbf{w}[h_i; b_i; h_i \odot b_i; h_i \odot \hat{b}_i] \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^{4d}$  is a trainable parameter;  $h_i \in H$  is the headline representation;  $\hat{b}_i$  represents the headline-to-article vector;  $[\cdot]$  stands for a vector concatenation across rows; and  $\odot$  represents element-wise multiplications.

Since contextual embeddings contain only high-level information, we believe that introducing pre-extracted human semantic knowledge into the computation of similarities will enhance the ability of the model to identify the boundaries of the clusters and is helpful to the final predictions. Therefore, we apply the pre-extracted human semantic knowledge in order to enhance context embeddings.

Next, the headline-aware body representation  $P$  is encoded by the neural network and a matrix  $U \in \mathbb{R}^{d \times n}$  is obtained as coarse-grained memories for words of the linked articles. The matrix  $U$  captures the interaction between the words of an articles with respect to its linked headline and is further processed by the knowledge-enhanced self-attention mechanism.

### 2) KNOWLEDGE-ENHANCED SELF-ATTENTION

Self-attention is a part of the fine-grained memory layer. Its purpose is to fuse the coarse-grained memories  $U$  into themselves. Through the knowledge-enhanced joint-attention layer, the headline-aware representation  $u_i^B$  of an article  $b_i \in B$  is generated to locate important parts of the linked article. However, such a representation has very limited knowledge of the context. To overcome this problem, we directly match the headline-aware representation of article  $b_i$  against itself [61], [62]. This process can effectively pinpoint important word pairs within the representation and update the representation with this useful information. The fine-grained memories  $v_i^B \in V$  of article  $b_i$  is computed as follows, where  $V \in \mathbb{R}^{d \times n}$  is a matrix that represents fine-grained memories of  $B$  and  $H$ .

$$v_i^B = \text{BiGRU}(v_{i-1}^B, [u_i^B, a_i]) \quad (2)$$

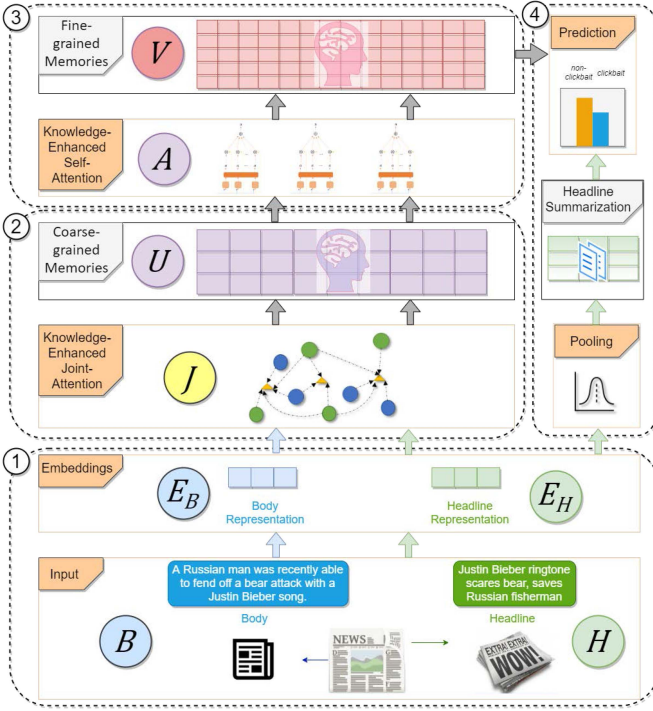


FIGURE 2. Our proposed end-to-end neural model.

where  $a_i = att(u^B, u_i^B)$  is an attention-pooling vector of the whole article ( $u^{B_i}$ ) and computed as follows:

$$a_i = \sum_{k=1}^n c_k^i u_k^B, \quad a_i \in A, \quad (3)$$

where

$$c_k^i = \frac{\exp(r_k^i)}{\sum_{j=1}^n \exp(r_j^i)} \quad (4)$$

$$r_j^i = p^\top \tanh(W_u^B u_j^B + W_u^{\hat{B}} u_i^B) \quad (5)$$

We apply an additional gate to  $[u_i^B, a_i]$  to adaptively control the input to the neural network [61]. Thus, a matrix  $V \in \mathbb{R}^{d \times n}$  is obtained as fine-grained memories.

Based on the relationships between words in the headline and its linked article, knowledge-enhanced self-attention extracts semantic information from the whole linked article.

#### IV. APPLYING THE PROPOSED ATTENTION-BASED NEURAL NETWORK TO CLICKBAIT DETECTION

As depicted in Fig. 2, our model contains four different layers to capture different linguistic representations. Following is a detailed description of the four layers.

- The first layer extracts information from headlines and their linked articles at the word level. We compute representations of words in the headlines and linked articles as follows [63]. For word embeddings, we use pretrained word vectors GloVe [56] for both headlines and their linked articles. In addition, we use the following two types of linguistic features for each token in the linked

articles: (i) named-entity recognizer (NER) encoding for 18 different types of the NER tags; (ii) part-of-speech (POS) tagging encoding for 56 different types of POS tags.

- The purpose of the second layer, a *coarse-grained memory layer*, is to apply both article-to-headline and headline-to-article attentions, based on representations of headlines and their linked articles. We generate preliminary memories over headline-article word pairs. Specifically, we apply *knowledge-enhanced joint-attention* (discussed in section III-B1) to fuse the body representation into the headline representation. Then we process this output with a BiGRU layer. Finally, the coarse-grained memories are obtained by concatenating the forward GRU outputs and the backward GRU outputs, which are the headline-aware body representation.
- In the third layer, a *fine-grained memory layer*, we construct the refined memories over headline-article word pairs using the coarse-grained memories. First, we apply *knowledge-enhanced self-attention* (discussed in section III-B2) to fuse coarse-grained memories into themselves. Then we process this output with a BiGRU layer. Finally, the fine-grained memories are obtained by concatenating the forward GRU outputs and the backward GRU outputs, which are the final body representation.
- The last layer, a *label prediction layer*, yields the label prediction based on the fine-grained memories and the headline context embeddings. Specifically, we apply attention pooling to the headline representation to obtain a summary of the headline. Finally, we compute a posterior distribution of the two candidate labels to determine if the headline is clickbait or not.

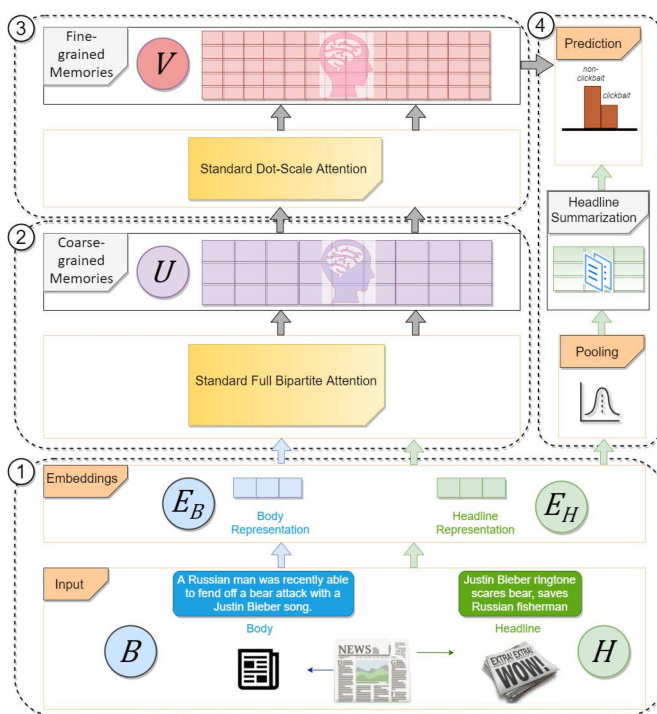
#### A. EXAMPLE AND VISUALIZATION

We present an example to qualitatively demonstrate the effectiveness of the proposed knowledge-enhanced attention mechanisms. We first implemented a base model by applying full bipartite attention (for joint attention) and standard dot-scale attention (for self attention) to the model, similar to the work in [5], [14]. These attention mechanisms in the base model do not utilize human semantic knowledge. As a result, we expected that it would not perform as well as the KED model that has knowledge-enhanced attention mechanisms. We name this base model NKED (No Knowledge Enhancement). The structure of NKED is depicted in Fig. 3. We then replaced full bipartite attention and standard dot-scale attention with the proposed knowledge-enhanced joint-attention and knowledge-enhanced self-attention mechanisms, respectively, to obtain the proposed model KED.

Table 2 shows an example of a headline and its linked article, and the visualization of the output from the knowledge-enhanced attention schemes. Qualitatively, we display the contribution of the phrases to the classification in terms of knowledge-enhanced self-attention. The most important phrases are highlighted in red with the intensity of the color indicating the degree of contribution. Meanwhile, we use waves to indicate the key phrases with high attention scores in terms

**TABLE 2.** Visualization of a Sample From the FNC Challenge Dataset. The Proposed KED Model is Implemented With Knowledge Enhancement to the Attention Mechanisms and NKED is Without Knowledge Enhancement. The Most Important Phrases in Self-Attention are Highlighted in Red Where the Intensity of the Color Indicates the Degree of Contribution. The Waves Indicate the Key Phrases With High Attention Scores in Knowledge-Enhanced Joint-Attention

Model		Sentence Samples	Prediction
NKED	Headline	Justin Bieber ringtone scares bear, saves Russian fisherman	clickbait ✘
	Body	A Russian man was recently able to fend off a bear attack with a Justin Bieber song.	
KED	Headline	Justin Bieber ringtone scares bear, saves Russian fisherman	non-clickbait ✔
	Body	A Russian man was recently able to fend off a bear attack with a Justin Bieber song.	



**FIGURE 3.** The diagram of NKED model.

of knowledge-enhanced joint-attention. (Detailed visualization techniques can be found in [64].)

As shown in Table 2, NKED (with no knowledge enhancement to the attention mechanisms) captures key phrases *saves* and *able to fend*, which have very low similarity. This mismatch contributes to NKED misclassifying the example as *clickbait*. In contrast, KED is able to correctly classify it as *non-clickbait*. The knowledge-enhanced joint-attention scheme first captures the global information described by the words *ringtone*, *song*, *scares*, *fend off*, *man* and *fisherman*. Informed by this global information, the knowledge-enhanced self-attention scheme reduces its attention to *saves*, while capturing key phrases {*Justin Bieber ringtone*, *Justin Bieber song*}, {*a Russian man*, *Russian fisherman*}, and {*scares bear*, *fend off a bear*}. In this case, the model infers that the

headline matches its linked article. Accordingly, the model makes a correct prediction labeled as *non-clickbait*.

In short, the global representation learned by the knowledge-enhanced joint-attention scheme provides an overall grasp of the whole text, which includes both semantic and structural information. It effectively helps the knowledge-enhanced self-attention scheme capture better instance-specific local features and improves classification performance.

In summary, the proposed neural network model is different from existing clickbait detection models in that it uses i) human semantic knowledge, which is pre-extracted using WordNet; ii) knowledge-enhanced attention mechanisms that enable the model to extract the most meaningful headline-article word pairs. On the one hand, the coarse-grained memory layer uses human semantic knowledge to assist both article-to-headline and headline-to-article attention. On the other hand, the fine-grained memory layer uses human semantic knowledge to assist the self-attention mechanism. As a result, the proposed model performs significantly better than existing state-of-the-art models, as will be shown next.

## V. EVALUATIONS

We present experimental results obtained from real-world datasets and compare the performance of the proposed model with that of existing state-of-the-art models/systems.

This section is organized as follows. Section V-A describes experiments settings. In Section V-B, we validate the effectiveness of pretrained word embeddings, knowledge enhancement to the attention mechanisms, and human semantic knowledge when applied to KED. Section V-C presents experimental results, comparing the performance of the proposed model with that of existing state-of-the-art systems and analyzing the results in detail. Section V-D verifies that when only a limited amount of training data is available (i.e., data scarcity), KED performs significantly better than the state-of-the-art clickbait detection models. Finally, Section V-E demonstrates the application of the proposed human semantic knowledge extraction scheme to several pretrained models in



**TABLE 3. Statistics of the Datasets.**

Dataset	Clickbait	Non-Clickbait	Clickbait
Clickbait Challenge	5,523	16,474	33.5%
FNC Challenge	54,894	20,491	72.8%

**TABLE 4. Model Settings: Vocabulary Size and Maximum Sequence Length**

Corpus	Vocabulary	Length	$\tau$
Clickbait Challenge	150,000	256	4
FNC Challenge	200,000	300	5

the semisupervised domain, and their significantly improved performance thanks to human semantic knowledge.

### A. EXPERIMENT SETTINGS

This section describes the datasets used in the experiments, the parameters of the proposed model, and the baseline models for comparison.

#### 1) DATASETS

We use two public annotated datasets for the experiments in this article.

*Clickbait Challenge*: This benchmark dataset [7] was released in 2017. It contains a total of 21,997 annotated samples, among which 17,598 are used for training and 4,399 for testing. Each sample is a pair of a headline and a linked article from Twitter posts. Five human annotators assigned a score to each sample, which ranges from 0 to 1, where ‘1’ denotes a post “heavily click-baiting” and ‘0’ indicates “not click-baiting” [65]. The samples with a mean score above 0.5 are considered as clickbait in this work [14].

*FNC Challenge*: This benchmark dataset was released by the Fake News Challenge competition in 2017 [3]. It contains 49,972 pairs of headline and linked article for training and 25,413 pairs for testing. Each pair of headline and linked article is classified into one of the four groups: ‘agree’, ‘discuss’, ‘disagree’, and ‘unrelated’. The samples with the label ‘unrelated’ are regarded as clickbait in this work [14].

Table 3 summarizes the statistics of the above datasets.

#### 2) MODEL SETTINGS

To implement the proposed neural network model, we use the Stanford CoreNLP [66] to preprocess the datasets. We apply the WordNet interface provided by NLTK [67] to extract human semantic knowledge. Additionally, we implement the proposed model using TensorFlow [68]. In the stage of extraction of human semantic knowledge, we set the hyperparameter  $\tau$  to 4 and 5 for the Clickbait Challenge and FNC Challenge dataset, respectively. We limit the vocabulary size to 150,000 and 200,000, and set the maximum sequence length to 256 and 300 for the Clickbait Challenge and FNC Challenge dataset, respectively. Table 4 summarizes the above parameters.

Following are the parameters common to both datasets. For the dense layers and the BiGRUs, we set the dimension  $d$  to

256. With respect to the training, we use Adam [69] as our optimizer. The learning rate was set to 0.001 and the mini-batch size, to 32. To avoid overfitting, we apply dropout [70] with a value of 0.1, and apply early stopping with a patience of 3. To avoid the exploding gradient problem, we apply gradient clipping [71] with a cutoff threshold of 2. We apply exponential moving averages with a decay rate of 0.999 to optimize the performance.

The above parameters were chosen to find the optimal settings for our experiments. Each data point in the reported results was averaged over three runs.

#### 3) BASELINE MODELS

We compare the proposed KED model with the following baseline models, which can be divided into the five following categories:

- *Deep feedforward networks*: DSSM [17]
- *CNN*: CLSM [18] and CBCNN [19]
- *RNN*: LiNN [4]
- *Attention-based networks*: MSA [13], BiGRU-ATT [5] and LSDA [14]
- *Transformer-based models*: LSACD [20], Transfer Learning [21]

The above five categories are listed in the first five blocks of Table 5. Following are brief descriptions of the baseline models.

*Deep Semantic Similarity Model (DSSM)*: Huang et al. [17] used deep neural networks to obtain hidden features of inputs and quantify the similarity in the space of latent representations. The authors preprocessed textual characteristics with n-gram schemes, and used the estimated similarities for prediction.

*Convolutional Latent Semantic Model (CLSM)*: Similar to DSSM, Shen et al. [18] leveraged convolutional neural networks for extracting inherent features.

*Linguistically-infused Neural Network (LiNN)*: Glenski et al. [4] introduced a model based on CNN and LSTM, which uses knowledge from both headlines and their linked articles. LSTM and CNN networks benefit the model by learning vectorized textual and visual knowledge individually.

*BiGRU-ATT*: Zhou [5] proposed an attention-based BiGRU model. The model first discovers inherent features of headlines and their linked articles. The discovered latent features are then consolidated to be fed into attention-based recurrent networks for the classification task.

*Multi-Strategy Approach (MSA)*: Kumar et al. [13] introduced a composite form of clickbait detection models. First, the authors applied an attention-based bidirectional RNN-based method to learn the input data. Then they incorporated the latent information with the relationship knowledge into Siamese networks for the final predictions.

*Clickbait Convolutional Neural Network (CBCNN)*: Zheng et al. [19] proposed a convolutional model that leverages only features from the headlines to detect clickbait. The authors first converted headlines into text vectors and then made predictions with textCNN [72].



**TABLE 5.** Performance Comparison Using the Clickbait Challenge and FNC Challenge Datasets. The Best Result in Each Column is Highlighted in green; Our Previously Published Best Results in [10], From  $KED_R$ , are Shown in blue; Other Researchers' Previously Published Best Results are Shown in red. The Results in the First to Fourth Block (From DSSM to LSDA) Come from [14]. The Results in the Fifth Block are Taken from [20], [21]. NKED: No Knowledge Enhancement to Attention Mechanisms.  $KED_R$ : With Knowledge Enhancement and Randomly Initialized Word Embeddings. KED: With Knowledge Enhancement and Pretrained GloVe Word Embeddings

Model	Clickbait Challenge				FNC Challenge			
	Accuracy	Precision	Recall	F1-score	Accuracy	Precision	Recall	F1-score
DSSM [17]	0.817	0.655	0.661	0.658	0.747	0.894	0.740	0.811
CLSM [18]	0.833	0.683	0.643	0.662	0.756	0.959	0.762	0.853
CBCNN [19]	0.844	0.654	0.653	0.653	0.789	0.852	0.845	0.857
LiNN [4]	0.827	0.642	0.621	0.631	0.868	0.925	0.884	0.913
MSA [13]	0.826	0.699	0.474	0.565	0.859	0.920	0.877	0.907
BiGRU-ATT [5]	0.856	0.719	0.650	0.683	0.879	0.924	0.897	0.919
LSDA [14]	0.860	0.722	0.699	0.710	0.894	0.933	0.912	0.928
LSACD [20]	0.880	-	-	0.740	-	-	-	-
Transfer Learning [21]	0.858	0.713	0.668	0.690	-	-	-	-
NKED	0.862	0.724	0.704	0.714	0.892	0.930	0.911	0.920
$KED_R$ [10]	0.880	0.754	0.756	0.755	0.913	0.948	0.935	0.941
KED	0.892	0.761	0.763	0.762	0.928	0.952	0.943	0.947

*Similarity-Aware Deep Attentive Model (LSDA)*: Dong et al. [14] introduced a deep learning model that is similarity-aware and attention-based to acquire and represent similarities.

*Lure and Similarity for Adaptive Clickbait Detection (LSACD)*: Zheng et al. [20] proposed a RNN model that examines the similarity between a headline and its linked content, and the headline's degree of enticement ("lure") to infer whether the headline is clickbait or not.

*Transfer Learning*: Rajapaksha et al. [21] adapted transformer models BERT, RoBERTa and XLNet for the clickbait detection task by applying several fine-tuning methods and model configuration changes.

The last three blocks of Table 5 list the following models that we implemented:

- NKED: randomly initialized word embeddings; without knowledge enhancement to the attention mechanisms.
- $KED_R$  [10]: randomly initialized word embeddings; with knowledge enhancement to the attention mechanisms.
- KED: pretrained word embeddings; with knowledge enhancement to the attention mechanisms.

## B. VALIDATIONS

In this section, we validate the effectiveness of pretrained word embeddings, knowledge enhancement to the attention mechanisms, and human semantic knowledge when applied to KED.

### 1) EFFECTIVENESS OF PRETRAINED WORD EMBEDDINGS

We investigate the effectiveness of pretrained word embeddings on KED. In our previous work [10], the model was implemented using randomly initialized word embeddings,

which we name  $KED_R$ . For this article, we replace randomly initialized word embeddings with pretrained GloVe vectors [56], expecting the latter will improve the classification performance of the model. We name the enhanced version KED. It is expected that GloVe can enhance the performance of KED over  $KED_R$ .

In this experiment, we use GloVe vectors with 300 dimensions to initialize the word embeddings. The results in terms of F1-score, accuracy, precision, and recall on the two datasets Clickbait Challenge and FNC Challenge are listed in Table 5.

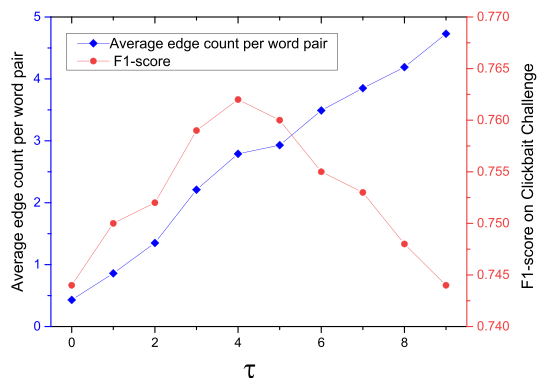
Overall, the performance of the proposed architecture can be noticeably improved by using pretrained GloVe word embeddings. For example, KED achieves a higher F1 score of 0.762, 0.7 percentage points (0.762 vs. 0.755) higher than  $KED_R$  on the Clickbait Challenge dataset, and of 0.947, 0.6 percentage points (0.947 vs. 0.941) higher on the FNC Challenge dataset. KED also performs much better than the state-of-the-art model LSDA. KED outperforms LSDA by 5.2 percentage points (0.762 vs. 0.710) and 1.9 percentage points (0.947 vs. 0.928) on the two datasets, respectively. More importantly, the performance gain is particularly significant on the relatively difficult Clickbait Challenge dataset. The experimental results confirm that pretrained word embeddings can improve the performance of our KED model.

### 2) EFFECTIVENESS OF KNOWLEDGE ENHANCEMENT TO ATTENTION MECHANISMS (ABLATION STUDY)

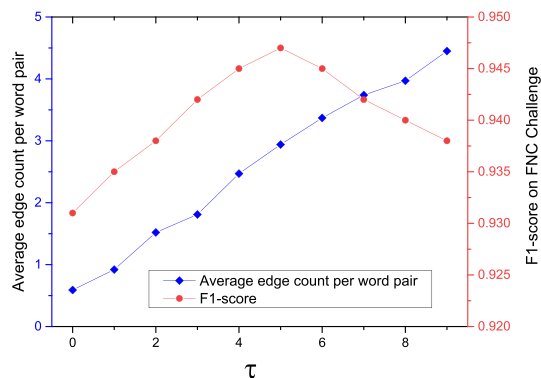
We conducted an ablation study by comparing the performance of  $KED_R$  (with knowledge enhancement to the attention mechanisms) to that of NKED (no knowledge enhancement). As depicted in Fig. 3, full bipartite attention and standard dot-scale attention are implemented in NKED.

**TABLE 6. Effectiveness of Human Semantic Knowledge. The Hyper-Parameter  $\tau$  Selected for the Best Result is Marked in Bold**

$\tau$	Clickbait Challenge		FNC Challenge	
	Average edge count per word pair	F1-score	Average edge count per word pair	F1-score
0	0.43	0.744	0.59	0.931
1	0.86	0.75	0.92	0.935
2	1.35	0.752	1.52	0.938
3	2.21	0.759	1.81	0.942
4	<b>2.79</b>	<b>0.762</b>	2.47	0.945
5	2.93	0.760	<b>2.94</b>	<b>0.947</b>
6	3.49	0.755	3.37	0.945
7	3.85	0.753	3.74	0.942
8	4.19	0.748	3.97	0.940
9	4.73	0.744	4.45	0.938



(a) Clickbait Challenge



(b) FNC Challenge

**FIGURE 4. Illustrations of the results given in Table 6.**

They are replaced by the proposed knowledge-enhanced joint-attention and knowledge-enhanced self-attention, respectively, in  $KED_R$ . Both  $KED_R$  and  $NKED$  use randomly initialized word embeddings.

The results are given in Table 5. We observe that  $KED_R$  performs significantly better than  $NKED$  in all cases, thanks to the knowledge enhancement. In particular, on the relatively difficult task Clickbait Challenge,  $KED_R$  outperforms  $NKED$  by 5.2 percentage points (0.756 vs. 0.704) in terms of recall, and by 4.1 percentage points (0.755 vs. 0.714) in terms of F1-score. The results demonstrate that human semantic knowledge is beneficial to the clickbait detection task, and plays an important role in our neural model.

### 3) EFFECTIVENESS OF HUMAN SEMANTIC KNOWLEDGE

To verify the effectiveness of applying human semantic knowledge to clickbait detection, we obtain ten augmented datasets by varying  $\tau$  from 0 to 9, and train a different  $KED$  system on each augmented clickbait detection dataset. Table 6 and Fig. 4 show that when the hyper-parameter  $\tau$  is increased from 0 to 9, the amount of human semantic knowledge (represented by the average edge count per word pair)

rises monotonically. The performance in terms of F1-score of  $KED$  rises until  $\tau$  reaches 4 and 5 for Clickbait Challenge and FNC Challenge, respectively. The reason is that a larger  $\tau$  value allows more hypernym relations to be extracted from the knowledge base as explained earlier in Section III-A. As  $\tau$  increases beyond a peak value (4 or 5), the performance of the model degrades. The reason is that very large values of  $\tau$  increase the chances for trivial hypernym relations to be extracted from the knowledge base, leading to lower performance. Thus it can be seen that human semantic knowledge provided by the WordNet-based human semantic knowledge extraction scheme plays an essential role in the training of the proposed  $KED$  system.

### C. PERFORMANCE COMPARISON WITH EXISTING WORKS

We compare the performance of the proposed  $KED$  model with that of the baseline models listed in Section V-A3 using the two benchmark datasets of clickbait detection. Table 5 summarizes the experimental results with four commonly used metrics: accuracy, recall, precision, and F1-score. The experimental results show that the proposed model  $KED$ , thanks to the use of human semantic knowledge and

knowledge-enhanced attention mechanisms, significantly outperforms the above state-of-the-art models for both clickbait detection benchmark datasets. For example, on the FNC Challenge dataset, KED outperforms LSDA by 3.4 percentage points in terms of accuracy (0.928 vs. 0.894) and by 3.1 percentage points in terms of recall (0.943 vs. 0.912), with LSDA having the previously published best results on this dataset.

It is worth noting that even the very basic version NKED (with randomly initialized word embeddings and without knowledge enhancement to attention mechanisms) performs better than most of the above state-of-the-art models. For example, on the relatively difficult task Clickbait Challenge, NKED outperforms the best performer of the fourth block (LSDA) by 0.5 percentage points (0.704 vs. 0.699) in terms of recall, and 0.04 percentage points (0.714 vs. 0.710) in terms of F1-score. The reason for the higher performance of NKED is its use of human semantic knowledge for the task of classification, while the previous models do not use human semantic knowledge.

We observe that the CNN and RNN based models in the second and third block perform better than the deep feed-forward networks in the first block. We believe the reason for this is the capability of CNN and RNN in capturing location information.

The models shown in the fourth block, MSA, BiGRU-ATT, and LSDA, perform better than those listed in the first three blocks thanks to the attention mechanisms, which can capture key phrases more effectively. The transformer-based models shown in the fifth block, LSACD and Transfer Learning, perform better than those listed in the first four blocks thanks to the transformer framework that capture key phrases even better. However, all the models underperform KED because they do not use human semantic knowledge or knowledge-enhanced attention mechanisms.

Human semantic knowledge combined with the proposed knowledge-enhanced attention mechanisms enables KED to outperform existing clickbait detection models. The following example illustrates the above advantages of KED over existing models. Consider the following two text samples:

- 1) "Little John was looking for his toy box. Finally, he found it. The box was in the pen."
- 2) "The first ballpoint model was patented by the American leather tanner and inventor John J. Loud."

There are no semantic correlations between 'pen' and 'ballpoint' in this context. The knowledge-enhanced attention mechanisms allow KED to assign a very low similarity score to the word pair {'pen,' 'ballpoint'} in this example, while previous models may (incorrectly) give it high similarity scores, leading to lower classification performance.

#### D. QUANTITATIVE ANALYSIS OF DATA SCARCITY

While deep learning does not require labour-intensive, time-consuming feature engineering compared to traditional machine learning techniques such as SVM, it usually requires a large amount of training data to work effectively. One of

the challenges researchers in deep learning often face is the lack of training data. Therefore, it is beneficial to know how a model performs when training data is scarce. In this experiment, we compare the performance of KED with that of state-of-the-art models in the case of data scarcity.

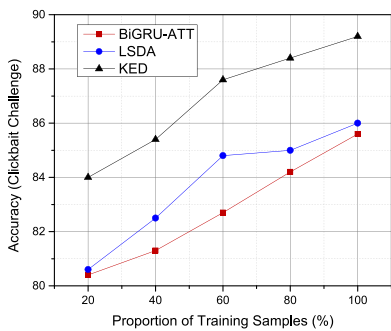
Instead of using all the training examples, we produced several subsets of training data of different sizes so as to study the relationship between the amount of available training data and the model performance. Specifically, from each of the two training datasets (Clickbait Challenge and FNC Challenge) we created four training subsets, which contain 20%, 40%, 60%, and 80% of the original training data set, respectively. We randomly selected headlines (and their linked articles) until the number of headlines reached 20% (40%, 60%, or 80%) of the size of an original data set (listed in Table 3). The random selections of headlines maintain the same clickbait and non-clickbait ratio of each original data set as listed in Table 3.

We compare KED with BiGRU-ATT and LSDA, and evaluate their performance in terms of F1-score, accuracy, precision, and recall using the subsets of training data mentioned above. We choose these two models to compare with KED because LSDA is among the previous top two performers and ranked first in most metrics; BiGRU-ATT is the first ranked model in the Clickbait Challenge competition [7]. Furthermore, BiGRU-ATT, which uses self-attention, and LSDA, which uses global similarity, are the most similar to our model.

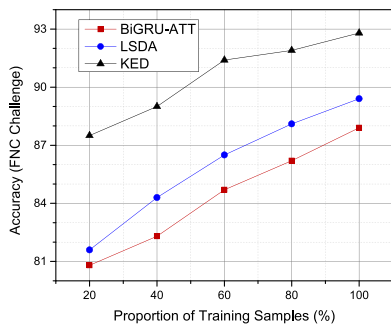
As shown in Fig. 5(a) through (h), as the percentage of training samples increases, the F1-score, accuracy, precision, and recall of all three models increase as expected. In all cases, KED performs much better than the other two. For instance, when 80% training samples are available, in terms of F1-score on the Clickbait Challenge dataset, KED achieves 0.738, 6.8 percentage points higher than LSDA (0.738 vs. 0.670) and 7.6 percentage points higher than BiGRU-ATT (0.738 vs. 0.662). Similarly, when 60% training samples are available, in terms of F1-score on the FNC Challenge dataset, KED achieves 0.920, 5.5 percentage points higher than LSDA (0.920 vs. 0.865) and 7.3 percentage points higher than BiGRU-ATT (0.920 vs. 0.847). In addition, when 80% training samples are available, in terms of recall on the Clickbait Challenge dataset, KED achieves 0.733, 5.3 percentage points higher than LSDA (0.733 vs. 0.680) and 9.1 percentage points higher than BiGRU-ATT (0.733 vs. 0.642). Overall, from the graphs, we can observe that our KED performs significantly better than the state-of-the-art clickbait detection models even when only a limited amount of training data is available.

#### E. INCORPORATING HUMAN SEMANTIC KNOWLEDGE INTO PRETRAINED MODELS IN THE SEMISUPERVISED DOMAIN

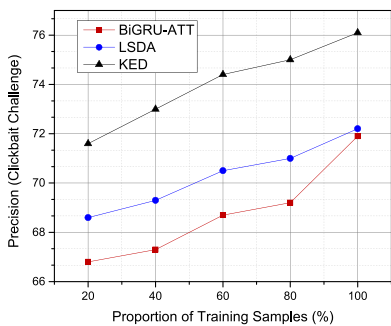
Our model KED is a fully supervised model, in which all model parameters are trained from scratch. On the other hand, several transformer-based pretrained NLP models in the semisupervised domain such as BERT [46], RoBERTa [47],



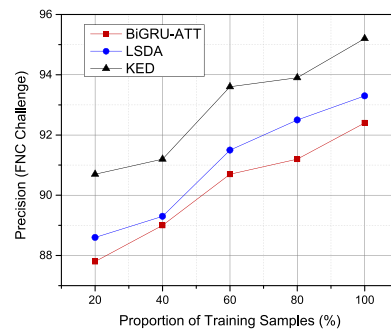
(a) Clickbait Challenge, accuracy



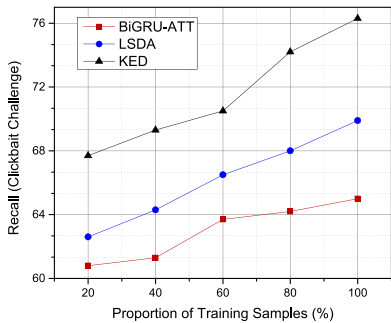
(b) FNC Challenge, accuracy



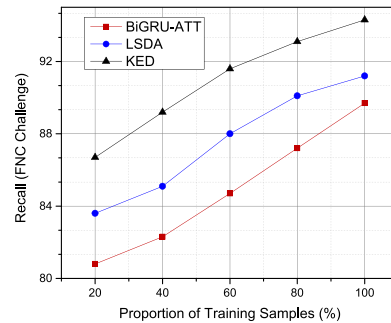
(c) Clickbait Challenge, precision



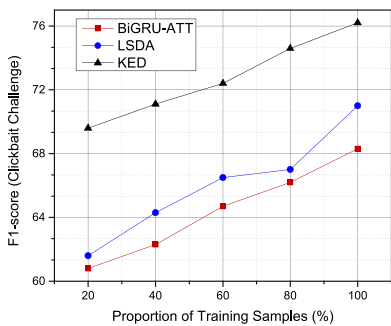
(d) FNC Challenge, precision



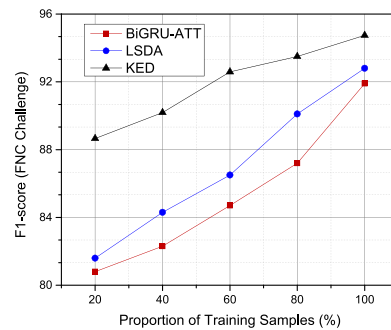
(e) Clickbait Challenge, recall



(f) FNC Challenge, recall



(g) Clickbait Challenge, F1-score



(h) FNC Challenge, F1-score

FIGURE 5. Illustrations of the quantitative analysis of data scarcity.



**TABLE 7.** Incorporation of Human Semantic Knowledge Into Pretrained Models BERT, RoBERTa, and XLNet

Model	Clickbait Challenge	FNC Challenge	Model	Clickbait Challenge	FNC Challenge
KED	0.892	0.928	KED	0.761	0.952
BERT (✗ ✓)	0.865   0.878	0.896   0.906	BERT (✗ ✓)	0.711   0.740	0.918   0.943
RoBERTa (✗ ✓)	0.871   0.882	0.902   0.913	RoBERTa (✗ ✓)	0.716   0.747	0.922   0.949
XLNet (✗ ✓)	0.877   0.895	0.906   0.929	XLNet (✗ ✓)	0.721   0.762	0.925   0.955
(a) Accuracy			(b) Precision		
Model	Clickbait Challenge	FNC Challenge	Model	Clickbait Challenge	FNC Challenge
KED	0.763	0.943	KED	0.762	0.947
BERT (✗ ✓)	0.683   0.713	0.896   0.919	BERT (✗ ✓)	0.697   0.726	0.907   0.931
RoBERTa (✗ ✓)	0.690   0.722	0.900   0.925	RoBERTa (✗ ✓)	0.703   0.734	0.911   0.937
XLNet (✗ ✓)	0.725   0.765	0.907   0.945	XLNet (✗ ✓)	0.723   0.763	0.916   0.950
(c) Recall			(d) F1-score		

The symbols “✓” and “✗” denote whether the proposed human semantic knowledge extraction scheme is incorporated into a model or not, respectively. The best result in each column is highlighted in orange. The second best result in each column is highlighted in green.

XLNet [48], ELECTRA [49] and ALBERT [50] have shown to be effective in learning common language representations and performing many NLP tasks [73]. Advantages of pre-trained NLP models over fully supervised models include quick and simple implementation of a classification model with acceptably good performance; much less labelled data required; and diverse use cases and applications. Therefore, they have been widely used in many NLP tasks [74]. We believe that human semantic knowledge can enhance the performance of pretrained models in the semisupervised domain. Therefore, we incorporated human semantic knowledge into BERT, RoBERTa, and XLNet, the most well-known pretrained NLP models, and conducted experiments to verify the effectiveness of human semantic knowledge on these models.

In order to incorporate human semantic knowledge, we use the proposed WordNet-based human semantic knowledge extraction scheme to generate additional annotated data to expand the original datasets Clickbait Challenge and FNC Challenge, and then fine-tune these pretrained models using the newly expanded datasets. We perform grid searches over the learning rate, epochs, and batch size in  $[2 \cdot 10^{-5}, 3 \cdot 10^{-5}, 5 \cdot 10^{-5}]$ ,  $[2, 3, 4]$  and  $[16, 32]$ , respectively. Additionally, the AdamW optimizer [75] is applied, and the weight decay and warm-up step are set to 0.01 and 0.05, respectively. Other than the above parameters, we re-use the hyper-parameters of the original models for both the original and enhanced (with human semantic knowledge) versions [46], [47], [48]. The results from the original pretrained models and the enhanced versions in terms of F1-score, accuracy, precision, and recall on the Clickbait Challenge and FNC Challenge datasets are listed in Table 7.

We note that after incorporating the proposed human semantic knowledge extraction scheme into the pretrained models BERT, RoBERTa, and XLNet, their performance is improved noticeably. On the Clickbait Challenge dataset, the F1 scores of BERT, RoBERTa, and XLNet are improved by 2.9 percentage points (0.726 vs. 0.697), 3.1 percentage points (0.734 vs. 0.703), and 4.0 percentage points (0.763 vs.

0.723), respectively. On the FNC Challenge dataset, the F1 scores of BERT, RoBERTa, and XLNet are improved by 2.4 percentage points (0.931 vs. 0.907), 2.6 percentage points (0.937 vs. 0.911), and 3.4 percentage points (0.950 vs. 0.916), respectively.

Additionally, after being enhanced with human semantic knowledge, all the pretrained models outperform the state-of-the-art model LSDA with both datasets and all metrics. For example, in terms of precision, BERT, RoBERTa, and XLNet outperform LSDA by 1.8 percentage points (0.740 vs. 0.722), 2.5 percentage points (0.747 vs. 0.722), and 4.0 percentage points (0.762 vs. 0.722), respectively, on the Clickbait Challenge dataset. Similarly, on the FNC Challenge dataset, BERT, RoBERTa, and XLNet outperform LSDA by 1.0 percentage points (0.943 vs. 0.933), 1.6 percentage points (0.949 vs. 0.933), and 2.2 percentage points (0.955 vs. 0.933), respectively. (LSDA uses fewer parameters than the pretrained models, however, tens of millions vs. hundreds of millions).

We also note that the proposed model KED with pretrained GloVe vectors outperforms the enhanced versions of BERT and RoBERTa, as shown in Table 7. The pretrained model XLNet with human semantic knowledge performs slightly better than KED, e.g., 0.3 percentage points higher (0.950 vs. 0.947), in terms of F1-score on the FNC Challenge dataset. However, XLNet is much more resource consuming than KED (and so are BERT and RoBERTa). For example, the pretrained models have more than 300 million parameters (335 million for BERT, 356 million for RoBERTa, 360 million for XLNet), whereas KED has much fewer parameters, in the range of tens of millions. In the next sub-section, we show that the training time and inference time of KED is significantly shorter than those incurred by BERT, RoBERTa and XLNet.

Overall, the experimental results in Table 7 show that the proposed human semantic knowledge extraction scheme significantly improves the performance of pretrained models in the semisupervised domain. Moreover, the proposed scheme is compatible with and easily applied to the pretrained models.

**TABLE 8. Training Time (In Minutes) on the Two Training Datasets**

	KED	BERT	RoBERTa	XLNet
Clickbait Challenge	15 min	22 min	25 min	29 min
FNC Challenge	28 min	36 min	38 min	41 min

**TABLE 9. Average Inference Time (In Milliseconds) of One Sample**

	KED	BERT	RoBERTa	XLNet
Clickbait Challenge	32 ms	51 ms	73 ms	88 ms
FNC Challenge	24 ms	40 ms	62 ms	75 ms

## F. EXECUTION TIME

As deep learning models become more complex, training and inference time is increased as more layers and more neurons per layer are added. Therefore, execution time should be considered to ensure that it is feasible to deploy a deep learning model for real life applications. In this sub-section, we report the training and inference time of the proposed KED model and the pretrained models BERT, RoBERTa and XLNet that are adapted for the clickbait detection task.

This set of experiments is the same as the set presented above in Section V-E. The execution time metrics are training time and inference time. All the experiments were run on a Windows workstation with the following configuration: NVIDIA GeForce 840 M graphics card, Intel Core i7-4710 2.5 GHz processor, 12 GB DDR3 memory, and 1 TB solid state drive.

Table 8 shows the training time of KED, BERT, RoBERTa and XLNet on the Clickbait Challenge and FNC Challenge datasets. Overall, we note that KED outperforms BERT, RoBERTa, and XLNet on both datasets. For example, on the Clickbait Challenge dataset, KED trained faster than BERT, RoBERTa, and XLNet by 47% (15 m vs. 22 m), 67% (15 m vs. 25 m), and 93% (15 m vs. 29 m), respectively.

To measure the inference time, we ran each model to classify 4,400 samples from the Clickbait Challenge dataset and 15,077 samples from the FNC Challenge dataset. The total inference time for classifying 4,400 samples from the Clickbait Challenge dataset was then divided by 4,400 to obtain the average inference time of one sample. The same calculation was done for the 15,077 samples from the FNC Challenge dataset. Table 9 lists the average inference time of one sample incurred by KED, BERT, RoBERTa and XLNet. We observe that the average inference time of KED is significantly shorter than that of BERT, RoBERTa, and XLNet on both datasets. For instance, on the Clickbait Challenge dataset, KED outperforms BERT, RoBERTa, and XLNet by 59% (32 vs. 51), 128% (32 vs. 73), and 175% (32 vs. 88), respectively.

It is worth noting that the classification performance of KED is only slightly lower than that of XLNet (accuracy of 0.928 vs. 0.929, and F1-score of 0.947 vs. 0.950 from 7), while the average inference time of KED is much shorter than that of XLNet, 2.8 to 3.1 times shorter (32 vs. 88, and 24 vs. 75, respectively, from 9). A similar comparison can also be

said for inference time. The much shorter execution time of KED is expected because it is designed and optimized for a specific task, while the pretrained models are intended to be adapted and used for a wide range of applications. It is also expected that transformer-based clickbait detection models such as LSACD [20], Transfer Learning [21] and LGBM [45], which adapted and fine tuned pretrained models, would incur longer execution time than KED.

The above results demonstrate that KED is feasible to be deployed in real world scenarios for the clickbait detection task thanks to reasonably short execution time. It provides the best of both worlds: classification performance (e.g., accuracy and F1-scores) and execution time.

## VI. CONCLUSION

We propose a neural network that extracts human semantic knowledge to build linguistic knowledge graphs which guide the attention mechanisms. The model can be used for many applications, one of which being clickbait detection for the purpose of combating fake news and illegitimate marketing. To the best of our knowledge, our work is the first that incorporates human semantic knowledge into the task of clickbait detection. We carried out extensive comparative experiments to evaluate the effectiveness and performance of the proposed attention-based neural network model. Experimental results show that

- the proposed KED model significantly outperforms state-of-the-art models/systems such as BiGRU-ATT, LSDA and LSACD;
- KED performs better than the state-of-the-art models even when training data is limited;
- pretrained word embeddings such as GloVe significantly improve the performance of KED compared to randomly initialized word embeddings;
- KED performs better or comparably to powerful pretrained models, namely, BERT, RoBERTa, and XLNet, while being much more lightweight and incurring significantly less training time and inference time.
- the use of human semantic knowledge significantly enhances not only the performance of KED in a fully supervised domain but also that of pretrained models in a semisupervised domain, namely, BERT, RoBERTa, and XLNet.

The proposed neural attention-based architecture can be relatively easily adapted to a variety of other NLP tasks such as word sense disambiguation and machine reading comprehension. In the future, we will study these applications using the above neural attention-based architecture. Currently all training data used with KED must be labelled. We will investigate an upgraded or new model that can work with a mix of labelled and unlabelled data [76] while maintaining all the advantages of KED.

## REFERENCES

- [1] Y. Chen, N. J. Conroy, and V. L. Rubin, "Misleading online content: Recognizing clickbait as 'false news,'" in *Proc. ACM Workshop Multimodal Deception Detection*, 2015, pp. 15–19.
- [2] M. Pothast, S. Köpsel, B. Stein, and M. Hagen, "Clickbait detection," in *Proc. Eur. Conf. Inf. Retrieval*, 2016, pp. 810–817.
- [3] "FNC challenge," 2017. [Online]. Available: <http://www.fakenewschallenge.org>
- [4] M. Glenski, E. Ayton, D. Arendt, and S. Volkova, "Fishing for clickbaits in social images and texts with linguistically-infused neural network models," 2017, *arXiv:1710.06390*.
- [5] Y. Zhou, "Clickbait detection in tweets using self-attentive network," 2017, *arXiv:1710.05364*.
- [6] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [7] "Clickbait challenge," 2017. [Online]. Available: <https://www.clickbait-challenge.org/>
- [8] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, and E. Hovy, "Hierarchical attention networks for document classification," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics: Hum. Lang. Technol.*, 2016, pp. 1480–1489.
- [9] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [10] F. Wei and U. T. Nguyen, "A neural attentive model using human semantic knowledge for clickbait detection," in *Proc. Int. Conf. Parallel Distrib. Process. Appl., Big Data Cloud Comput., Sustain. Comput. Commun., Social Comput. Netw.*, 2020, pp. 770–776.
- [11] Q. Chen, X. Zhu, Z.-H. Ling, D. Inkpen, and S. Wei, "Neural natural language inference models enhanced with external knowledge," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguistics*, 2018, pp. 2406–2417.
- [12] P. Biyani, K. Tsioutsoulouklis, and J. Blackmer, "8 amazing secrets for getting more clicks: Detecting clickbaits in news streams using article informality," in *Proc. 13th AAAI Conf. Artif. Intell.*, 2016, pp. 94–100.
- [13] V. Kumar, D. Khattar, S. Gairola, Y. Kumar Lal, and V. Varma, "Identifying clickbait: A multi-strategy approach using neural networks," in *Proc. 41st Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2018, pp. 1225–1228.
- [14] M. Dong, L. Yao, X. Wang, B. Benatallah, and C. Huang, "Similarity-aware deep attentive model for clickbait detection," in *Proc. Pacific-Asia Conf. Knowl. Discov. Data Mining*, 2019, pp. 56–69.
- [15] L. Chen, M. A. L. Ralph, and T. T. Rogers, "A unified model of human semantic knowledge and its disorders," *Nature Hum. Behav.*, vol. 1, no. 3, pp. 1–10, 2017.
- [16] G. A. Miller, "WordNet: A lexical database for English," *Commun. ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [17] P.-S. Huang, X. He, J. Gao, L. Deng, A. Acero, and L. Heck, "Learning deep structured semantic models for web search using clickthrough data," in *Proc. 22nd ACM Int. Conf. Inf. Knowl. Manage.*, 2013, pp. 2333–2338.
- [18] Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil, "A latent semantic model with convolutional-pooling structure for information retrieval," in *Proc. 23rd ACM Int. Conf. Conf. Inf. Knowl. Manage.*, 2014, pp. 101–110.
- [19] H.-T. Zheng, J.-Y. Chen, X. Yao, A. K. Sangaiah, Y. Jiang, and C.-Z. Zhao, "Clickbait convolutional neural network," *Symmetry*, vol. 10, no. 5, 2018, Art. no. 138.
- [20] J. Zheng, K. Yu, and X. Wu, "A deep model based on lure and similarity for adaptive clickbait detection," *Knowl.-Based Syst.*, vol. 214, 2021, Art. no. 106714.
- [21] P. Rajapaksha, R. Farahbakhsh, and N. Crespi, "BERT, XLNet or RoBERTa: The best transfer learning model to detect clickbaits," *IEEE Access*, vol. 9, pp. 154704–154716, 2021.
- [22] H. Allcott and M. Gentzkow, "Social media and fake news in the 2016 election," *J. Econ. Perspectives*, vol. 31, no. 2, pp. 211–36, May 2017.
- [23] J. Aro, "The cyberspace war: Propaganda and trolling as warfare tools," *Eur. View*, vol. 15, no. 1, pp. 121–132, 2016.
- [24] D. J. Watts, D. M. Rothschild, and M. Mobius, "Measuring the news and its impact on democracy," *Proc. Nat. Acad. Sci.*, vol. 118, no. 15, 2021, Art. no. e1912443118.
- [25] C. Carvalho, N. Klagge, and E. Moench, "The persistent effects of a false news shock," *J. Empirical Finance*, vol. 18, no. 4, pp. 597–615, 2011.
- [26] S. Kogan, T. J. Moskowitz, and M. Niessner, "Fake news: Evidence from financial markets," 2019, doi: [10.2139/ssrn.3237763](https://doi.org/10.2139/ssrn.3237763).
- [27] W.-Y. Sylvia Chou, A. Gaysynsky, and J. N. Cappella, "Where we go from here: Health misinformation on social media," *Amer. J. Public Health*, vol. 110, no. S3, pp. S273–S275, 2020.
- [28] H. Pan, "A study on deep learning: Training, models and applications" Ph.D. dissertation, Dept. Electr. Eng. Comput. Sci., York Univ., Toronto, ON, Canada, 2017.
- [29] W. S. McCulloch and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *Bull. Mathe. Biophys.*, vol. 5, no. 4, pp. 115–133, 1943.
- [30] F. Rosenblatt, "The perceptron: A probabilistic model for information storage and organization in the brain," *Psychol. Rev.*, vol. 65, no. 6, 1958, Art. no. 386.
- [31] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [32] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Netw.*, vol. 4, no. 2, pp. 251–257, 1991.
- [33] Y. Zhu et al., "Aligning books and movies: Towards story-like visual explanations by watching movies and reading books," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 19–27.
- [34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [35] X. Wang et al., "Truth discovery via exploiting implications from multi-source data," in *Proc. 25th ACM Int. Conf. Inf. Knowl. Manage.*, 2016, pp. 861–870.
- [36] P. Bourgonje, J. M. Schneider, and G. Rehm, "From clickbait to fake news detection: An approach based on detecting the stance of headlines to articles," in *Proc. EMNLP Workshop: Natural Lang. Process. Meets Journalism*, 2017, pp. 84–89.
- [37] A. Chakraborty, B. Paranjape, S. Kakarla, and N. Ganguly, "Stop clickbait: Detecting and preventing clickbaits in online news media," in *Proc. IEEE/ACM Int. Conf. Adv. Social Netw. Anal. Mining*, 2016, pp. 9–16.
- [38] H.-T. Zheng, X. Yao, Y. Jiang, S.-T. Xia, and X. Xiao, "Boost clickbait detection based on user behavior analysis," in *Proc. Asia-Pacific Web, Web-Age Inf. Manage., Joint Conf. Web Big Data*, 2017, pp. 73–80.
- [39] B. Vijgen et al., "The Listicle: An exploring research on an interesting shareable new media phenomenon," *Studia Universitatis Babeş-Bolyai-Ephemerides*, vol. 59, no. 1, pp. 103–122, 2014.
- [40] A. Anand, T. Chakraborty, and N. Park, "We used neural networks to detect clickbaits: You won't believe what happened next!," in *Proc. Eur. Conf. Inf. Retrieval*, 2017, pp. 541–547.
- [41] M. M. U. Rony, N. Hassan, and M. Yousuf, "Diving deep into clickbaits: Who use them to what extents in which topics with what effects?," in *Proc. IEEE/ACM Int. Conf. Adv. Social Netw. Anal. Mining 2017*, pp. 232–239.
- [42] A. Razaque, B. Alotaibi, M. Alotaibi, S. Hussain, A. Alotaibi, and V. Jotsov, "Clickbait detection using deep recurrent neural network," *Appl. Sci.*, vol. 12, no. 1, 2022, Art. no. 504.
- [43] A. Razaque et al., "Blockchain-enabled deep recurrent neural network model for clickbait detection," *IEEE Access*, vol. 10, pp. 3144–3163, 2022.
- [44] M. Zhou, W. Xu, W. Zhang, and Q. Jiang, "Leverage knowledge graph and GCN for fine-grained-level Clickbait detection," *World Wide Web*, vol. 25, no. 3, pp. 1243–1258, 2022.
- [45] M. Marreddy, S. R. Oota, L. S. Vakada, V. C. Chinni, and R. Mamidi, "Clickbait detection in Telugu: Overcoming NLP challenges in resource-poor languages using benchmarked techniques," in *Proc. Int. Joint Conf. Neural Netw.*, 2021, pp. 1–8.
- [46] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics: Hum. Lang. Technol.*, 2019, pp. 4171–4186.
- [47] Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," 2019, *arXiv:1907.11692*.
- [48] Z. Yang et al., "XLNet: Generalized autoregressive pretraining for language understanding," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 5754–5764.
- [49] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "ELECTRA: Pre-training text encoders as discriminators rather than generators," in *Proc. Int. Conf. Learn. Representations*, 2020, pp. 1–18.
- [50] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A lite BERT for self-supervised learning of language representations," in *Proc. Int. Conf. Learn. Representations*, 2020, pp. 1–17.



- [51] R. Rehurek and P. Sojka, "Software framework for topic modelling with large corpora," in *Proc. LREC Workshop New Challenges NLP Frameworks*, 2010, pp. 45–50.
- [52] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 3111–3119.
- [53] M. Dong, L. Yao, X. Wang, B. Benatallah, Q. Z. Sheng, and H. Huang, "DUAL: A deep unified attention model with latent relation representations for fake news detection," in *Proc. Int. Conf. Web Inf. Syst. Eng.*, 2018, pp. 199–209.
- [54] M. Coleman and T. L. Liau, "A computer readability formula designed for machine scoring," *J. Appl. Psychol.*, vol. 60, no. 2, 1975, Art. no. 283.
- [55] J. Anderson, "Lix and Rix: Variations on a little-known readability index," *J. Reading*, vol. 26, no. 6, pp. 490–496, 1983.
- [56] J. Pennington, R. Socher, and C. Manning, "GloVe: Global vectors for word representation," in *Proc. Empirical Methods Natural Lang. Process.*, 2014, pp. 1532–1543.
- [57] M. Seo, A. Kembhavi, A. Farhadi, and H. Hajishirzi, "Bidirectional attention flow for machine comprehension," 2016, *arXiv:1611.01603*.
- [58] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998–6008.
- [59] A. W. Yu et al., "QANet: Combining local convolution with global self-attention for reading comprehension," in *Proc. Int. Conf. Learn. Representations*, 2018, pp. 1–16.
- [60] W. Chen, X. Quan, C. Kit, Z. Min, and J. Wang, "Multi-choice relational reasoning for machine reading comprehension," in *Proc. 28th Int. Conf. Comput. Linguistics*, 2020, pp. 6448–6458.
- [61] W. Wang, N. Yang, F. Wei, B. Chang, and M. Zhou, "Gated self-matching networks for reading comprehension and question answering," in *Proc. 55th Annu. Meeting Assoc. Comput. Linguistics*, 2017, pp. 189–198.
- [62] T. Rocktäschel, E. Grefenstette, K. M. Hermann, T. Kočiský, and P. Blunsom, "Reasoning about entailment with neural attention," in *Proc. Int. Conf. Learn. Representations*, 2016, pp. 1–9.
- [63] X. Liu, Y. Shen, K. Duh, and J. Gao, "Stochastic answer networks for machine reading comprehension," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguistics*, 2018, pp. 1694–1704.
- [64] J. Li, X. Chen, E. Hovy, and D. Jurafsky, "Visualizing and understanding neural models in NLP," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics: Hum. Lang. Technol.*, 2016, pp. 681–691.
- [65] M. Potthast, T. Gollub, M. Hagen, and B. Stein, "The clickbait challenge 2017: Towards a regression model for clickbait strength," 2018, *arXiv:1812.10847*.
- [66] C. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. Bethard, and D. McClosky, "The stanford CoreNLP natural language processing toolkit," in *Proc. 52nd Annu. Meeting Assoc. Comput. Linguistics: Syst. Demonstrations*, 2014, pp. 55–60.
- [67] S. Bird and E. Loper, "NLTK: The natural language toolkit," in *Proc. ACL Interactive Poster; Demonstration Sessions*, 2004, pp. 214–217.
- [68] M. Abadi et al., "TensorFlow: A system for large-scale machine learning," in *Proc. 12th USENIX Symp. Operating Syst. Des. Implementation*, 2016, pp. 265–283.
- [69] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.
- [70] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [71] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *Proc. Int. Conf. Mach. Learn.*, 2013, pp. 1310–1318.
- [72] K. Yoon, "Convolutional neural networks for sentence classification," in *Empirical Methods Natural Lang. Process.*, 2014, pp. 1746–1751.
- [73] F. Wei and U. T. Nguyen, "Stock trend prediction using financial market news and BERT," in *Proc. 12th Int. Conf. Knowl. Discov. Inf. Retrieval*, 2018, pp. 325–332.
- [74] X. Han et al., "Pre-trained models: Past, present and future," *AI Open*, vol. 2, pp. 225–250, 2021.
- [75] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *Proc. Int. Conf. Learn. Representations*, 2019, pp. 1–18.
- [76] K. Chen, L. Yao, D. Zhang, X. Wang, X. Chang, and F. Nie, "A semi-supervised recurrent convolutional attention model for human activity recognition," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 5, pp. 1747–1756, May 2020.



**FENG WEI** (Student Member, IEEE) is currently working toward the Ph.D. degree with the Department of Electrical Engineering and Computer Science, York University, Toronto, ON, Canada. His research interests include machine learning, cybersecurity, natural language processing, and computer vision.



**UYEN TRANG NGUYEN** (Member, IEEE) received the bachelor degree in computer science and the master's degree in computer science from Concordia University, Montreal, QC, Canada, and the Doctoral degree in computer science from the University of Toronto, Toronto, ON, Canada. She is currently an Associate Professor with the Lassonde School of Engineering, York University, Toronto, ON, Canada. Her research interests include wireless networking, mobile computing, online social networking, information security, and financial technology. Her research has been funded by NSERC National Sciences and Engineering Research Council of Canada, MITACS Mathematics of Information Technology and Complex Systems, and industry partners. She is a frequent reviewer of grant applications for NSERC, MITACS, the U.S. National Science Foundation, and the National Science Centre of Poland.