

Harnessing the Power of Text-image Contrastive Models for Automatic Detection of Online Misinformation

Hao Chen¹, Peng Zheng¹, Xin Wang^{2*}, Shu Hu³, Bin Zhu⁴, Jinrong Hu^{1*}, Xi Wu¹, Siwei Lyu²

¹Chengdu University of Information Technology, China ²University at Buffalo, USA

³Carnegie Mellon University, USA ⁴Microsoft, USA

{haochen,hjr,wuxi}@cuit.edu.cn {xwang264,siweilyu}@buffalo.edu
3210607015@stu.cuit.edu.cn shuhu@cmu.edu binzhu@microsoft.com

Abstract

As growing usage of social media websites in the recent decades, the amount of news articles spreading online rapidly, resulting in an unprecedented scale of potentially fraudulent information. Although a plenty of studies have applied the supervised machine learning approaches to detect such content, the lack of gold standard training data has hindered the development. Analysing the single data format, either fake text description or fake image, is the mainstream direction for the current research. However, the misinformation in real-world scenario is commonly formed as a text-image pair where the news article/news title is described as text content, and usually followed by the related image. Given the strong ability of learning features without labelled data, contrastive learning, as a self-learning approach, has emerged and achieved success on the computer vision. In this paper, our goal is to explore the contrastive learning in the domain of misinformation identification. We developed a self-learning model and carried out the comprehensive experiments on a public data set named COSMOS. Comparing to the baseline classifier, our model shows the superior performance of non-matched image-text pair detection (approximately 10%) when the training data is insufficient. In addition, we observed the stability for contrastive learning and suggested the use of it offers large reductions in the number of training data, whilst maintaining comparable classification results.

1. Introduction

The proliferation of news articles on social media platforms allows for real-time access to information, but also leads to an increase in the spread of misinformation due to deceptive practices. Misinformation has an adverse impact



Figure 1. An example of text-image pairing. For the image on the left, the text in the green box is the **real** (matched) caption, while the text in the red box is a **fake** (random) caption.

on both cyber and physical societies, and has gained considerable coverage in the last several years. For example, the COVID-19 pandemic has provided a fertile ground for conspiracy theories. One of the most prominent ones was that the 5G technology was somehow responsible for the emergence of the novel coronavirus. This theory gained particular momentum in early April 2020 and led to a wave of vandalism targeting communication infrastructure in the UK and other nations [14]. In addition, misinformation causes anxiety in populations across different ages [38].

Therefore, there is an urgent demand for tools that can detect fake news content accurately and efficiently in order to eliminate misinformation and protect the harmony of online environment. Manually monitoring authenticity of information requires substantial human efforts and time [45]. In recent years, dominant automatic approaches [21–23, 27, 36, 45] rely on machine learning techniques, particularly supervised classification models, to identify fake content. However, labelled data sets that supervised learning needs are very limited in both sizes and diversities, which presents a major obstacle for applications of supervised learning algorithms in this domain. Determining whether a news article is fake or not is generally a chal-

*corresponding author

lenging task since it requires comprehensive knowledge and verification by authentic sources.

Additionally, most researchers concentrates on using natural language processing (NLP) techniques to identify fake text-based content, overlooking the fact that news articles often contain both text and images. An example is shown in Figure. 1. The caption in the green box is the real description of the image while the caption in the red box is false. Such misinformation is easily manipulated or created by an AI-based generator and then rapidly circulated through the Internet. Relying on textual semantic and syntactic similarities [1, 18, 35] can achieve promising identification performance to some extent. However, these models ignore relationships across multiple data modals, particularly the image-text, which can deteriorate the accuracy of identification. To address the limitations of these models, Aneja et al. [3] attempt to solve the problem by combing both data types. Specifically, their dataset, COSMOS, consists of images and two captions for each image, and their task is to predict if the two captions are both corresponding to the image, i.e., the out-of-context (OOC) classification. They use a convolutional neural network as an image encoder and a pre-trained language model as a textual caption encoder, and achieve 85% classification accuracy on COSMOS. However, the model is trained on the large corpus which is less efficient and time-consuming. In addition, results mostly are influenced by the pre-trained model (SBERT) according to our investigation.

In this paper, we extend Aneja et al.’s method [3] by using a language-vision model based on contrastive learning [8] for out-of-context detection on the COSMOS dataset. As a self-learning approach, the contrastive learning shows a strong ability to learn feature representations without annotating a large-scale dataset. It learns representations of data by contrasting similar and dissimilar pairs of examples. The merit of contrastive learning lies in its ability to leverage the inherent structure of data to learn more useful representations. By contrasting similar examples and pulling them closer together in representation space, and pushing dissimilar examples further apart, contrastive learning can learn more robust and discriminative features that can be used for a variety of downstream tasks, such as image classification or language understanding. This has been shown to be especially effective in settings where labeled data is scarce or expensive to obtain, as contrastive learning can learn from large amounts of unlabeled data. We use Aneja et al.’s [3] method as the baseline in our experimental evaluation. Our results indicate that our proposed method outperforms the baseline on the out-of-context prediction. The main contributions of this paper can be summarized as follows:

- Our study on the baseline model reveals for the first time that the baseline relies primarily on the textual

similarity from the pre-trained model (SBERT) to classify OOC, which may potentially introduce bias and distort the results;

- We present a new model incorporating a contrastive learning component, which we will show is advantageous in capturing image feature representations, particularly when the training data is limited in size;
- We conduct extensive experiments to evaluate our proposed method and compare with the baseline method. These experiments demonstrate that our method outperforms the baseline steadily in a varying training data sizes.
- We developed the classification model to identify the correct caption.

We note that this work focuses on the same prediction task as the baseline on the COSMOS dataset. It is not the real-world misinformation detection wherein the prediction is to predict if a pair of image and its caption is consistent or not. This is because we do not have a dataset that can mimic the real-world misinformation detection. That said, we should point out that our proposed method with minor modifications can be applied to the scenario of the real-world misinformation detection. Actually, most of our proposed method is general and can be applied to other classification tasks. In particular, our use of contrastive learning is beneficial for scenarios when there is a lack of labeled training data.

The paper is organized as follows: Section 2 reviews related work, and Section 3 describes our proposed method. Our experimental evaluation is presented in Section 4. The conclusion and future work are presented in Section 5.

2. Related Work

Online misinformation has become a topic of interests over the past few years, motivating the research community to address the problem. Bondielli et al. [4] categorise information as fake or rumours depending on whether the news has been confirmed by the authoritative sources. Both Guo et al. [20] and Meel et al. [32] elaborate the differences of various terms related to misinformation on social media, such as hoax, disinformation, and fake news. Instead of nitpicking on nuances of different definitions, we would like to focus on machine learning techniques themselves, specifically contrastive learning. Therefore, we decide to use ‘out-of-context’ to refer to misinformation comprising inconsistent image-text pairs. In this context, we briefly review related work, which includes misinformation detection and contrastive learning.

2.1. Misinformation Detection

Early research on the automatic identification of misinformation on social media concentrates on single-type data, particularly textual content. Traditional supervised clas-

sification techniques have been widely used in this area, such as support vector machine (SVM) [13, 17, 29], naïve bayes [40], logistic regression [9], and decision tree [6, 17]. Classic feature representations such as bag of words and n-grams with TF-IDF are generally used, with semantic or syntactic information ignored due to individually treated features (word tokens). This issue is subsequently alleviated by using feature engineering, a process of extracting and adding both linguistic and handcrafted features manually. For example, Pelrine et al. [35] systematically compare a set of transformer-based models on textual misinformation detection across various social media data sets. They point out the benefit of feature engineering. Shu et al. [40] investigate social networks and use spatiotemporal information such as numbers of retweets, timestamp, and locations to improve classifier results. Kwon et al. [30] explore users' profiles to increase the detection accuracy. However, feature engineering requires human efforts, in particular the knowledge of linguistic and social science. In addition, Zhu et al. [47] point out entities in news articles can change over time, which adversely impacts detection results. Inspired by the emerge of word2vec and paragraph2vec, a plenty of recent research [1, 18, 33, 39, 46] explore distributed representations where text content is converted into a dense vector by a language-embedding model, which is usually pre-trained on a general language corpus and thus preserves intrinsic language features such as syntactic and semantic features.

In addition to textual misinformation, another widely spread form of online news misinformation is de-contextualization (aka. out-of-context pairing) where images and their associated texts are unrelated to each other. Many researchers apply multi-modal analysis to tackle the problem of detecting this type of misinformation. Singhal et al. [42] introduce an ensemble model that exploits both textual and visual features to identify image-text fake news. Likewise, Giachanou et al. [16] combine image features extracted by a VGG [41] model and text features by a Bert [11] model to detect image-text misinformation. Recently, Aneja et al. [3] focus on the "cheapfake" generated by using AI-free approaches, such as filtering, re-contextualizing, and photo-shopping, rather than on "deepfake" generated by using AI-based techniques. They suggest utilizing image and text embeddings to forecast if an image-caption pair is out-of-context. They also release a substantial dataset for further research on this matter and it is becoming a well-known dataset in the media forensic domain. Nonetheless, their model was trained on the entire dataset. This paper aims to build upon their research by examining the efficiency of the contrastive learning model when training data is restricted.

2.2. Contrastive Learning

When using machine learning techniques for classification tasks, such as out-of-context (OOC) detection, data needs to be converted to a compact feature. Over the past decade, the dominant approach for determining image features is learning in a supervised way, such as training from ImageNet [10]. ViT [12] is a widely used visual feature representation model that uses the transformer [43] framework and is trained on classification tasks. Inspired by the achievement of BERT [11] in the NLP domain, computer vision community starts to increasingly focus on unsupervised training. Contrastive learning, as a self-learning approach, has gained popularity because it is able to learn feature representation without annotated data. Contrastive learning aims to move augmented samples generated from the same sample close to each other while keeping samples from different data far away.

Many contrastive learning models have been proposed [31]. He et al. [25] propose Momentum Contrast (MoCo) for unsupervised visual representation learning that matches encoded data with a serious of keys using the contrastive loss [34]. Subsequently, Chen et al. [8] propose SimCLR that generates training instances by separating different data augmentations. They comprehensively analyze a variety of image manipulation methods such as crop, resize, flip, color distort, rotate, cutout, and Gaussian noise. Gao et al. [15] propose SimCSE by adapting SimCLR to textual data that generates positive instances by different Bert dropouts and takes other samples within the batch as negative instances. Caron et al. [5] propose SwAV (Swapping Assignments between Views) by modifying SimCLR that clusters data and leverages contrastive learning techniques without requiring the computation of individual augmented samples. Radford et al. [37] introduce CLIP (Contrastive Language-Image Pre-training) that connects an image with a text description and creates a feature space for both data types. The model provides an efficient way to generate a text based on an image or vice versa. He et al. [24] create a masked autoencoders (MAE) to reconstruct masked patches when an image is split into multiple patches. Instead of training a model by using both positive and negative samples as aforementioned methods do, Grill et al. [19] propose BYOL that depends only on positive samples and is able to achieve outstanding performance for the feature representation.

3. Our Proposed Method

Our goal is to leverage the power of contrastive learning in the feature representation to identify inconsistent text-image pairing on social media. We describe its training and testing details in this section.

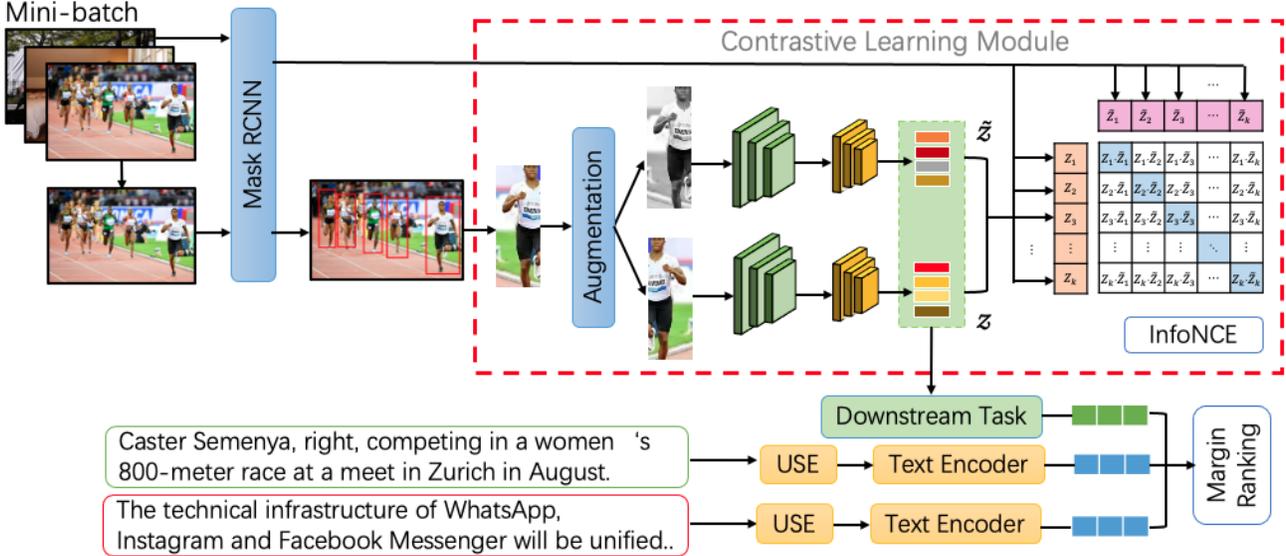


Figure 2. The structure of multi-modal contrastive learning. The model is accomplished by two branches, image feature learning and classification training. First, a picture in mini-batch (the size is 64) has been fed into a Mask RCNN by which the 10 objects have been detected. Each object is then augmented individually, followed by the full connected layer to produce the dense vector denoted as \tilde{z} and z . Subsequently, the matrix of the mini-batch is developed for the InfoNCE loss training, where the pair of z and \tilde{z} from the same object is treated as positive instances (the diagonal of the matrix), otherwise, the rest pairs are negative instances (see Eq. 1). In the classification training, the two captions – matched (green) and another caption sampled randomly (red) – are encoded using the Universal Sentence Encoder model (USE) [7]. The output of text encoder combining with the output of image encoder are passed to compute the similarities between object-caption pairs and finally used to reduce the margin ranking loss as following Eq. 3.1.

3.1. Contrastive Learning-guided Image-text Matching Training

Inspired by the baseline model on COSMOS [2], we extract features from images and texts separately and interact them to learn their matching. The training procedure is shown in Figure 2 and described in detail as follows.

Contrastive Learning Module. For each image, we use pre-trained Masked-RCNN [26] as the object detection [28] backbone to detect objects included in the image. Then we feed images and their detected objects (bounding boxes) into the augmentation module.

Within the augmentation module, each detected object is augmented. Augmented images are then fed to an image encoder followed by a full connected layer to generate a dense vector. We consider all augmented images from the same sample as positive instances and randomly selected images from other samples of the dataset as negative instances for training the contrastive learning model. Specifically, applying Mask RCNN on the input image, we can obtain N detected objects to form a set $\{x_k\}_{k=1}^N$ of objects, where $x_k \in \mathbb{R}^{d_x}$ and d_x is the dimension of detected object x . Then we apply data augmentation twice to get $2N$ objects, $\{\tilde{x}_l\}_{l=1}^{2N}$, where \tilde{x}_{2k} and \tilde{x}_{2k-1} are two random augmentations of x_k ($k = 1, \dots, N$). Different augmentation strategies can be used, such as rotations, adding noise, transla-

tion, brightness, etc. Thus an object can be augmented to generate more than two augmented objects, but only two augmented objects are used for each detected object in our setting. We use the following notation for two related augmented objects. Let $i \in I := \{1, \dots, 2N\}$ be the index of an arbitrary augmented object and $j(i)$ is the other augmented object that shares the same source object as the i -th augmented object. We feed augmented objects to an object encoder, represented as $E(\cdot) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^d$ (d is the output dimension of E), which is a ResNet-50 backbone followed by three components: RoIAlign, average pooling, and two fully-connected (FC) layers. Then, we can obtain a 300-dimensional vector for each augmented object, which maps the object feature representation into the application space of the contrastive loss, i.e., $\tilde{z}_i = E(\tilde{x}_i)$ and $\tilde{z}_{j(i)} = E(\tilde{x}_{j(i)})$ for two augmented objects from the same source object.

To shorten the distance between encoder vectors \tilde{z}_i and $\tilde{z}_{j(i)}$ from the same source object and widen the distance between \tilde{z}_i and an augmented object from another source object, we use the self-supervised contrastive learning to formulate the self-supervised contrastive loss \mathcal{L}_{CL} as follows,

$$\mathcal{L}_{CL} = \frac{-1}{|I|} \sum_{i \in I} \log \frac{\exp(\tilde{z}_i \cdot \tilde{z}_{j(i)} / \tau)}{\sum_{a \in A(i)} \exp(\tilde{z}_i \cdot \tilde{z}_a / \tau)}, \quad (1)$$

Algorithm 1 The Out-of-Context Matching

Data: sample $\{X_k\}_{k=1}^N$ in batch size N
 $X_k = \langle caption_r \rangle image \langle caption_m \rangle$
 τ and γ are constant

for all $k \in \{1, \dots, N\}$ **do**
 $O_m \leftarrow \text{MaskRCNN}(X_k)$ \triangleright 10 object detection
 $C_{km} \leftarrow t \cdot (X_{km})$ \triangleright text encoding for matched
 $C_{kr} \leftarrow t \cdot (X_{kr})$ \triangleright text encoding for random

for all $m \in \{1, \dots, 10\}$ **do**
 $\{A, \hat{A}\} \leftarrow A(O_m)$ \triangleright augmentation
 $Z \leftarrow f \cdot (A)$ \triangleright image encoding
 $\hat{Z} \leftarrow f \cdot (\hat{A})$

end for

end for

$M = N * 10$ \triangleright #of augmentation in batch

for all $i \in \{1, \dots, 2M\}$ and $j \in \{1, \dots, 2M\}$ **do**
 $s_{i,j} = z_i z_j / (\|z_i\| \|z_j\|)$ \triangleright pairwise similarity

end for

define $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2M} 1_{[k \neq i]} \exp(s_{i,k}/\tau)}$

$\mathcal{L}_{CL} = \frac{1}{2M} \sum_{k=1}^M [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$

define $s_m = \max(Z_k \cdot C_{km} | k \in \{1, \dots, N\})$

define $s_r = \max(Z_k \cdot C_{kr} | k \in \{1, \dots, N\})$

$\mathcal{L}_{Match} = [s_r - s_m + \gamma]_+$
 update networks f and t to minimize \mathcal{L}_{CL} and \mathcal{L}_{Match}

where $|I|$ is the cardinality of I , $\tau \in \mathbb{R}^+$ is a positive scalar temperature parameter, \cdot is the inner (dot) product operator, and $A(i) := I \setminus \{i\}$. It is common to regard i as an anchor. $j(i)$ is called the *positive* and the other $2N - 2$ indices ($\{k \in A(i) \setminus \{j(i)\}\}$) are called the *negatives*. The numerator in the log function of Eq. (1) is the representation distance between \tilde{z}_i and $\tilde{z}_{j(i)}$. The denominator is the representation distance between \tilde{z}_i and a total of $2N - 1$ terms, including the positive and negatives. With this contrastive learning module, we can enhance the accuracy of representations from the encoder.

Image-text Matching Module. This module match an image and its text (i.e., caption). Given matched caption c_m of the input image and a random caption c_r from a different image in the dataset, we follow [2] to use a pre-trained transformer-based Universal Sentence Encoder (USE, $U(\cdot)$) [7] to encode captions into unified 512-dimensional vectors. The vectors are then sent to an additional text encoder ($T(\cdot)$) to convert to a specific dimensional feature space \mathbb{R}^d that matches the the output dimension (i.e., d) of $E(\cdot)$. In particular, the text encoder is a ReLU followed by one FC layer. Therefore, we can represent $\tilde{c}_m = T(U(c_m))$ and $\tilde{c}_r = T(U(c_r))$ for the final embedded features of c_m and c_r , respectively.

Then we evaluate the match performance of the object embedding and the caption embedding. Specifically, we use

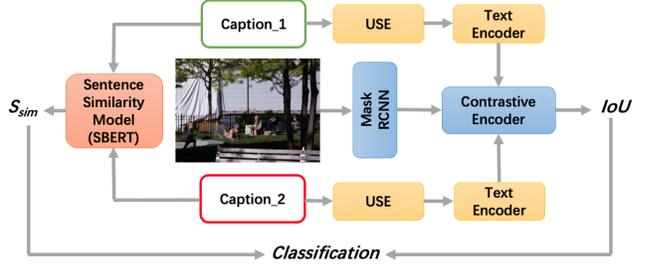


Figure 3. The testing structure. IoU indicates whether the two captions are describing the same object and S_{sim} represents the semantic similarity between the two captions. It predicts out-of-content if both scores are higher than their preset thresholds.

dot product to calculate the similarity between \tilde{z}_i and \tilde{c}_m (or \tilde{c}_r). We extract the maximum value as the final similarity score,

$$s_m = \max(\{\tilde{z}_i^T \tilde{c}_m | i \in I\}),$$

$$s_r = \max(\{\tilde{z}_i^T \tilde{c}_r | i \in I\}),$$

where s_m is the final similarity score for the matched caption and s_r is the final similarity score for the random caption. Our goal is to make s_m as larger as possible and s_r as smaller as possible. Thus, we design the following max-margin loss for training the model,

$$\mathcal{L}_{Match} = [s_r - s_m + \gamma]_+, \quad (2)$$

where $[a]_+ = \max(0, a)$ is the hinge function. $\gamma \in \mathbb{R}$ is a preset margin hyperparameter. The algorithm is showed in Algorithm 1.

3.2. Image-text Mismatching Training

Cross Training. We first train object encoder E in the contrastive learning module based on \mathcal{L}_{CL} (Eq. (1)) for all images in the dataset. Then we fix the contrastive learning module and train text encoder T according to \mathcal{L}_{Match} (Eq. (2)) on all images. The weights of the whole model are updated iteratively.

Joint Training. In addition to cross training, we explore joint training as well. Rather than freezing one of the loss functions during training, we normalize the loss of contrastive learning module \mathcal{L}_{CL} (Eq. (1)) and add \mathcal{L}_{Match} (Eq. (2)) to it to get the overall average loss on all images.

3.3. Image-text Mismatching Prediction

We follow [2] for our model's prediction of mismatching image and text or not, as shown in Figure 3. The prediction is determined by two scores, the IoU score and the Sentence BERT (SBERT) score (S_{sim}). The former score indicates whether the two captions are describing the same image region (object), and the latter is calculated for their sentence similarity.

Specifically, given a testing data (an image and two captions, e.g. $\langle caption_1 \rangle image \langle caption_2 \rangle$), we use the state-of-the-art SBERT model [44], which is pre-trained on the Sentence Textual Similarity (STS) task, to get the (S_{sim}) score for assessing both semantic and syntactic similarities between two the sentences. SBERT takes two pieces of text content as input and output a score in the range from 0 to 1. A higher score indicates that the two captions share more similar context.

For the IoU score, we use both image encoder and text encoder that are obtained from the trained language-vision model. We first compute the visual correspondences of objects B_{IC_i} in the image for each caption respectively. For example, B_{IC_1} represents the largest value (object) of image-caption alignment for I and $caption_1$. Then we calculate IoU for the overlapping of two bounding box (area) corresponding to $caption_1$ and $caption_2$. At last, we apply the following rule to predict the image-caption pair is out-of-context only if the two captions are describing the same object while they are having the different semantic meaning:

- Out-of-Content, if $IoU(B_{IC_1}, B_{IC_2}) > threshold$ and $S_{sim} < threshold$;
- Non Out-of-Context, otherwise.

4. Experimental Evaluation

4.1. Datasets & Pre-processing

As aforementioned, the lack of gold-standard training data is a major obstacle for us to investigate our approaches. To avoid this issue, this paper decided to use the data set that is large-scale and publicly available. Aneja et al. [3] create the data from two primary sources, fact-checking website and various mainstream news platforms such as New York Times, CNN, Reuters and so on. The original data is documented as JSON-format file. The structure of data is formed as $\langle caption_1 \rangle image \langle caption_2 \rangle$ where each image is followed by two captions, the one is genuine and other one is synthetic. The data summary is showing as Table 1.

We acknowledged that there is no labelling process in the training set and validation set as the synthetic caption for each image is randomly chose from the rest of text descriptions. However, the text-image pairs in testing set are manually labelled by the authors. An example of data instance is showing in Figure 1.

For the text pre-processing, we carried out the entity extraction to replace all the names, locations, and dates to the unique tokens respectively. For instance, the caption “Caster Semenya, right, competing in a women’s 800-meter race at a meet in Zurich in August.” is changed to “Caster GPE, right, competing in a women’s QUANTITY race at a meet in GPE in DATE.”

	#of Images	#of Captions	Annotation
Training Data	160k	360k	no
Validation Data	40k	90k	no
Testing Data	1700	1700	yes

Table 1. The statistic summary of data sets

4.2. Experiments Setup

This section will elaborate the experiments setup. Overall, we have conducted two experiments to demonstrate the benefit of contrastive learning when tackling small size of training set, and also a new model for identifying which of caption is true. To be specific, we compared three different models in our research:

- **baseline**, the model is originated from the paper [2].
- **cross training**, we use the contrastive learning structure to generate the feature representation for image as shown in Figure 2. Since the model have two loss functions (*InfoNCE* and *MarginRanking*), we iteratively freeze the one of them and train on the other.
- **joint training**, we further replaced the contrastive learning model to the joint training where the two loss functions (*InfoNCE* and *MarginRanking* are normalized and reduced at same time).

Evaluation Metric. Given the ultimate goal of this paper is to boost the ability of detecting OOC content, we use the standard classification evaluation metrics (*Accuracy*, *Precision*, *Recall* and *F1-score*).

Implementation Details. In addition, we acknowledge that detecting OOC is a trade-off issue where we have to decide whether the cost of false negatives (OOC has been classified as non-OOC) is higher than the false positives (non-OOC has been classified as OOC). We assume that in a real-world scenario, failing to identify misinformation and allowing it to spread would have a greater impact on social networks than incorrectly labeling clean content as misinformation. In this context, we will extend the original research that is only showing accuracy, and put more focus on the recall (also known as true positive rate).

Considering the neural network is randomly initialized, we carried out our experiments three times and get the averages. In addition, we used the following hyper-parameters as default in proposed model: Adam optimizer, ReLU activation function, a batch size of 64, and training for 10 epochs.

4.3. Contrastive Learning vs. Baseline

Rather than using the full size of training data, we would like to investigate the performance of aforementioned models on the small training set. We used the configurations

	Accuracy	Precision	Recall	F1-Score
<i>Baseline</i>	0.73	0.74	0.87	0.80
<i>Cross Training</i>	0.80	0.75	0.9	0.82
<i>Joint Training</i>	0.80	0.74	0.92	0.82

Table 2. The comparison results of three models over the Accuracy, Precision, Recall and F1-score.

as following, *Mask-RCNN*: 10 objects are detected; *Augmentation*: we chose one of noises such as rotation, adding gray, filtering, resizing, translation, brightness, etc.; *Resnet*: 18 convolutional layers are implemented and the output is 512 dimension (also text encoder is 512 dimension); *Dense Layer*: the output of dense vector is 300 dimension.

The results have been presented in Table 2, highlighting the best performance for each evaluation metric. It is observed that the baseline model performs the worst across all four metrics. As anticipated, replacing the structure of the original convolutional network with a contrastive learning module in the image encoder improves the ability to detect OOC content. Both cross training and joint training models achieve an accuracy of 0.80, which is an improvement of approximately 10% over the baseline model (0.73). In terms of recall, the contrastive learning models achieve better results than the baseline as well, with cross training and joint training models achieving 0.9 and 0.92 respectively. The increased performance is attributed to the enlargement of the training samples through augmentation, indicating that the contrastive learning models can effectively address the issue of insufficient training data. However, based on the paired t-test, the difference between cross training and joint training contrastive models is not statistically significant (p-value is 0.718).

We recognize that the accuracy of the aforementioned methods is lower than the results reported in the original paper, primarily because we randomly selected a subset of the training data, which consisted of less than 5,000 samples. In comparison to the reported accuracy of 85% achieved by using the full training set of 16,000 samples, we demonstrate that contrastive learning can achieve nearly 94% (0.80/0.85) of the performance using only 28% (4.5k/16k) of the training data. Furthermore, adding more data for training would yield limited improvements.

4.4. Contrastive Learning for True Caption Classification

We identified that most of classification decision is based on the sentence similarity BERT model which is pre-trained without fine-tuned. As shown in Figure 3, the final prediction is decided by the IoU and S_{sim} scores. However, the majority of testing data (more than 80%) has nearly 0.9 of IoU which is above the threshold (0.5). Consequently, the final classification is mainly attributes to the S_{sim} . To alle-

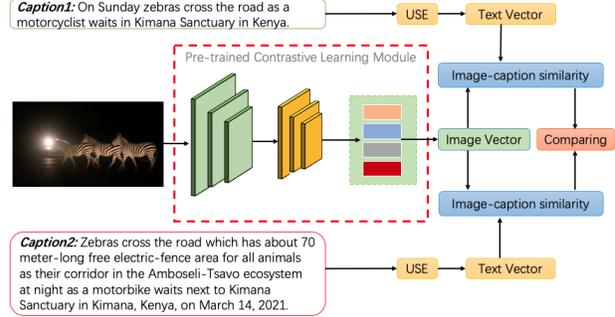


Figure 4. Classification of true caption based on the image-caption model. The pre-trained contrastive learning module (see Figure 2) is used for the image feature representation, and two captions are encoded through the Universal Sentence Encoder respectively. Then the caption vectors are used to compute the similarity with image vector. Finally, the prediction is decided by the similarity score where the image vector is closer to the true caption than the false one.

viate the bias, we would like to directly use contrastive in the training phase.

We simply modify our model for the task of classifying the correct caption. The structure is illustrated in Figure 4. Firstly, we extract the vector representation for the image from the trained contrastive learning model. Subsequently, we calculate the cosine similarity between this representation and the vector representations of the two captions (matched and non-matched captions). Finally, we determine the true caption based on the similarity score where the matched caption should have higher value with the image. The results obtained are mediocre, as only 941 out of 1700 (55%) have been correctly detected. We will continue working on improving this method in the future work.

4.5. Comparison on Varying Training Data Sizes

Following the results we obtained from previous experiment (Section 4.3), we would like further to investigate the impact of varying levels of training size on performance of the models. To be detailed, we want to examine whether the contractive learning would achieve comparable results when training on the even smaller size of data, and how the various mount of training data would consequent the models' classification ability. To do this, we created 10 different levels of training sizes, ranging from 500 to 5000 samples in intervals of 500 (randomly selected from original data set). For each level, we use the identical training data across three models. Then, we measured classification accuracy on the testing set as the same way with the previous experiment. The results are shown in Figure 6.

Theoretically, increasing the training data should improve the models' classification performance. However, according to the Figure 6, the baseline model showed the op-

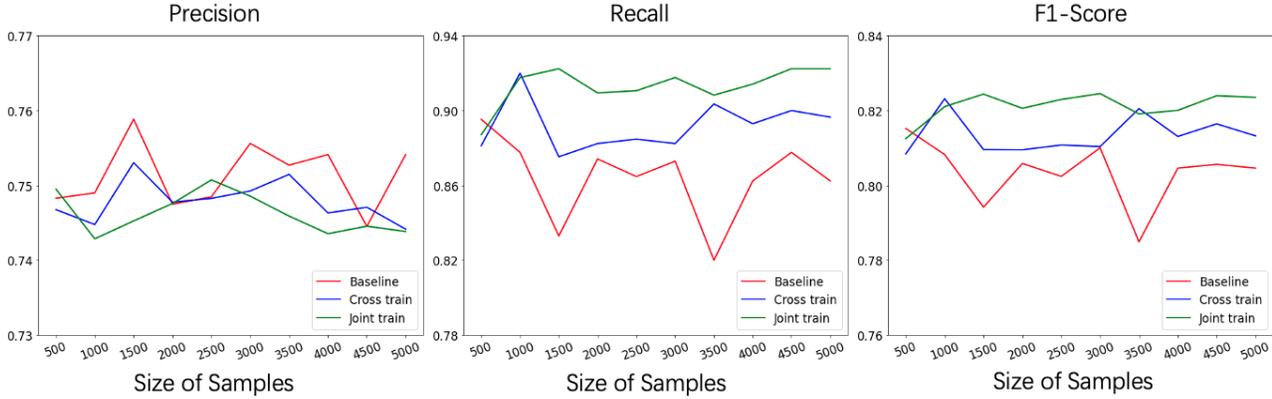


Figure 5. OOC classification precision, recall and F1-score for three different models with varying size of training samples.

posite trend. For instance, the best accuracy (0.78) of red line was achieved at 1500 training samples, and the baseline model is surprisingly decreased multiple times at training over the 2000, 3000, and 4500 samples. All three drops were significant, with the accuracy decreasing from nearly 0.78 to 0.72. Although the baseline models were fluctuating, the overall trend was upward (red line increased from 0.72 to 0.77).

In contrast, both contrastive learning models achieved superior results (around 0.79) when trained on very few data samples. For instance, the results of the two contrastive models fell in the range of 0.78 to 0.80, while the lowest accuracy of the baseline model was 0.72. Moreover, their performance steadily improved with the addition of more data, although the improvement was limited.

Overall, the two contrastive learning models performed comparably, and it was challenging to determine which one was better. To further explore this, we also reported the precision, recall, and f1-score in Figure 5. The joint training achieved the best results in terms of recall and f1-score, followed by the cross training, while the baseline model had the worst performance.

5. Conclusions

The focus of this work was to investigate the performance of contrastive learning for the feature representation when the tackling the domain of labelled data insufficiency, specifically the text-image context pairing. We highlight the following points from this work: (1) We proposed an advanced model for the task of out-of-context (OOC) detection based on the contrastive learning which is a self-supervised machine learning technique and utilize data augmentation for training. According to the results, we demonstrated the superior of our model comparing to the benchmark from original paper; (2) We focused on the situation where the labelled data is inadequate for training, which is a general limitation for the most of classification tasks.

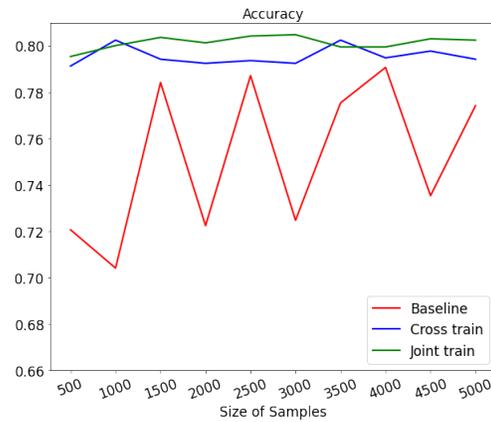


Figure 6. OOC classification accuracy for three different models with varying size of training samples.

We have done the comparison and show that the contrastive learning has a strong ability to learn the image feature and achieve the 94% of full performance, albeit decreasing the training data size nearly 70%; (3) We carried out a comprehensive analysis of the ability of the different classifiers to deal with the varying training data sizes. Our results show that the contrastive learning model produce the stable performance and increase the accuracy steadily once adding the training data. However, the baseline model has ups-and-downs results.

At last but not the least, we notice that one of the disadvantages for the COSMOS data set is that the label of OOC is based on whether the two captions are consistent with each other and also corresponding to the image. In this case, it ignores the real scenario of judging which of the caption is true. Although we proposed a new classifier model to deal with this issue, the results are not promising.

In the future, we would like further to explore contrastive learning and focus on the advance techniques to improve the misinformation detection accuracy.

References

- [1] Mohammed N Alenezi and Zainab M Alqenaei. Machine learning in detecting covid-19 misinformation on twitter. *Future Internet*, 13(10):244, 2021. 2, 3
- [2] Shivangi Aneja, Christoph Bregler, and Matthias Nießner. Cosmos: Catching out-of-context misinformation with self-supervised learning. 2021. 4, 5, 6
- [3] Shivangi Aneja, Cise Midoglu, Duc-Tien Dang-Nguyen, Sohail Ahmed Khan, Michael Riegler, Pål Halvorsen, Chris Bregler, and Balu Adsumilli. Acm multimedia grand challenge on detecting cheapfakes. *arXiv preprint arXiv:2207.14534*, 2022. 2, 3, 6
- [4] Alessandro Bondielli and Francesco Marcelloni. A survey on fake news and rumour detection techniques. *Information Sciences*, 497:38–55, 2019. 2
- [5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. 3
- [6] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web*, pages 675–684, 2011. 3
- [7] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. Universal Sentence Encoder. *arXiv e-prints*, page arXiv:1803.11175, Mar. 2018. 4, 5
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 2, 3
- [9] Alton YK Chua and Snehasish Banerjee. Linguistic predictors of rumor veracity on the internet. In *Proceedings of the International MultiConference of Engineers and Computer Scientists*, volume 1, page 387. Nanyang Technological University Singapore, 2016. 3
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3
- [13] Omar Enayet and Samhaa R El-Beltagy. Niletmrgr on semeval-2017 task 8: Determining rumour and veracity support for rumours on twitter. In *Proceedings of the 11th international workshop on semantic evaluation (SemEval-2017)*, pages 470–474, 2017. 3
- [14] Sarah Evanega, Mark Lynas, Jordan Adams, Karinne Smolenyak, and Cision Global Insights. Coronavirus misinformation: quantifying sources and themes in the covid-19 ‘infodemic’. *JMIR Preprints*, 19(10):2020, 2020. 1
- [15] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021. 3
- [16] Anastasia Giachanou, Guobiao Zhang, and Paolo Rosso. Multimodal multi-image fake news detection. In *2020 IEEE 7th International Conference on Data Science and Advanced Analytics (DSAA)*, pages 647–654. IEEE, 2020. 3
- [17] Georgios Giasemidis, Colin Singleton, Ioannis Agrafiotis, Jason RC Nurse, Alan Pilgrim, Chris Willis, and Danica Vukadinovic Greetham. Determining the veracity of rumours on twitter. In *International Conference on Social Informatics*, pages 185–205. Springer, 2016. 3
- [18] Sherry Girgis, Eslam Amer, and Mahmoud Gadallah. Deep learning algorithms for detecting fake news in online text. In *2018 13th international conference on computer engineering and systems (ICCES)*, pages 93–97. IEEE, 2018. 2, 3
- [19] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3
- [20] Bin Guo, Yasan Ding, Lina Yao, Yunji Liang, and Zhiwen Yu. The future of misinformation detection: new perspectives and trends. *arXiv preprint arXiv:1909.03654*, 2019. 2
- [21] Hui Guo, Shu Hu, Xin Wang, Ming-Ching Chang, and Siwei Lyu. Eyes tell all: Irregular pupil shapes reveal gan-generated faces. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2904–2908. IEEE, 2022. 1
- [22] Hui Guo, Shu Hu, Xin Wang, Ming-Ching Chang, and Siwei Lyu. Open-eye: An open platform to study human performance on identifying ai-synthesized faces. In *2022 IEEE 5th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 224–227. IEEE, 2022. 1
- [23] Hui Guo, Shu Hu, Xin Wang, Ming-Ching Chang, and Siwei Lyu. Robust attentive deep neural network for detecting gan-generated faces. *IEEE Access*, 10:32574–32583, 2022. 1
- [24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3
- [25] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 3
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. 4

- [27] Shu Hu, Yuezun Li, and Siwei Lyu. Exposing gan-generated faces using inconsistent corneal specular highlights. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2500–2504. IEEE, 2021. 1
- [28] Shu Hu, Chun-Hao Liu, Jayanta Dutta, Ming-Ching Chang, Siwei Lyu, and Naveen Ramakrishnan. Pseudoprop: Robust pseudo-label generation for semi-supervised object detection in autonomous driving systems. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4390–4398, 2022. 4
- [29] Shu Hu, Xin Wang, and Siwei Lyu. Rank-based decomposable losses in machine learning: A survey. *arXiv preprint arXiv:2207.08768*, 2022. 3
- [30] Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th international conference on data mining*, pages 1103–1108. IEEE, 2013. 3
- [31] Dizhong Lin, Ying Fu, Xin Wang, Shu Hu, Bin Zhu, Qi Song, Xi Wu, and Siwei Lyu. Contrastive class-specific encoding for few-shot object detection. In *2022 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2022. 3
- [32] Priyanka Meel and Dinesh Kumar Vishwakarma. Fake news, rumor, information pollution in social media and web: A contemporary survey of state-of-the-arts, challenges and opportunities. *Expert Systems with Applications*, 153:112986, 2020. 2
- [33] Jamal Abdul Nasir, Osama Subhani Khan, and Iraklis Varlamis. Fake news detection: A hybrid cnn-rnn based deep learning approach. *International Journal of Information Management Data Insights*, 1(1):100007, 2021. 3
- [34] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 3
- [35] Kellin Pelrine, Jacob Danovitch, and Reihaneh Rabbany. The surprising performance of simple baselines for misinformation detection. In *Proceedings of the Web Conference 2021*, pages 3432–3441, 2021. 2, 3
- [36] Wenbo Pu, Jing Hu, Xin Wang, Yuezun Li, Shu Hu, Bin Zhu, Rui Song, Qi Song, Xi Wu, and Siwei Lyu. Learning a deep dual-level network for robust deepfake detection. *Pattern Recognition*, 130:108832, 2022. 1
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 3
- [38] Yasmim Mendes Rocha, Gabriel Acácio de Moura, Gabriel Alves Desidério, Carlos Henrique de Oliveira, Francisco Dantas Lourenço, and Larissa Deadame de Figueiredo Nicolete. The impact of fake news on social media and its influence on health during the covid-19 pandemic: A systematic review. *Journal of Public Health*, pages 1–10, 2021. 1
- [39] Somya Ranjan Sahoo and Brij B Gupta. Multiple features based approach for automatic fake news detection on social networks using deep learning. *Applied Soft Computing*, 100:106983, 2021. 3
- [40] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context and spatialtemporal information for studying fake news on social media. *arXiv preprint arXiv:1809.01286*, 2018. 3
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- [42] Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin’ichi Satoh. Spofake: A multi-modal framework for fake news detection. In *2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM)*, pages 39–47, 2019. 3
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [44] Bin Wang and C.-C Kuo. Sbert-wk: A sentence embedding method by dissecting bert-based word models. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, PP:1–1, 07 2020. 6
- [45] Xin Wang, Hui Guo, Shu Hu, Ming-Ching Chang, and Siwei Lyu. Gan-generated faces detection: A survey and new perspectives. *arXiv preprint arXiv:2202.07145*, 2022. 1
- [46] Apurva Wani, Isha Joshi, Snehal Khandve, Vedangi Wagh, and Raviraj Joshi. Evaluating deep learning approaches for covid19 fake news detection. In *International Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situation*, pages 153–163. Springer, 2021. 3
- [47] Yongchun Zhu, Qiang Sheng, Juan Cao, Shuokai Li, Danding Wang, and Fuzhen Zhuang. Generalizing to the future: Mitigating entity bias in fake news detection. *arXiv preprint arXiv:2204.09484*, 2022. 3