

# Good Visual Guidance Makes A Better Extractor: Hierarchical Visual Prefix for Multimodal Entity and Relation Extraction

Xiang Chen<sup>1,2</sup>, Ningyu Zhang<sup>1,2\*</sup>, Lei Li<sup>1,2</sup>, Yunzhi Yao<sup>1,2</sup>, Shumin Deng<sup>1,2</sup>,  
Chuanqi Tan<sup>3</sup>, Fei Huang<sup>3</sup>, Luo Si<sup>3</sup>, Huajun Chen<sup>1,2,\*</sup>

<sup>1</sup>Zhejiang University & AZFT Joint Lab for Knowledge Engine, China

<sup>2</sup>Hangzhou Innovation Center, Zhejiang University, China

<sup>3</sup>Alibaba Group, China

{xiang\_chen, zhangningyu, leili21, yyztodd, 231sm, huajunsir}@zju.edu.cn,

{chuanqi.tcq, f.huang, luo.si}@alibaba-inc.com

## Abstract

Multimodal named entity recognition and relation extraction (MNER and MRE) is a fundamental and crucial branch in information extraction. However, existing approaches for MNER and MRE usually suffer from error sensitivity when irrelevant object images incorporated in texts. To deal with these issues, we propose a novel **Hierarchical Visual Prefix fusion NeTwork (HVPNeT)** for visual-enhanced entity and relation extraction, aiming to achieve more effective and robust performance. Specifically, we regard visual representation as plugable visual prefix to guide the textual representation for error insensitive forecasting decision. We further propose a dynamic gated aggregation strategy to achieve hierarchical multi-scaled visual features as visual prefix for fusion. Extensive experiments on three benchmark datasets demonstrate the effectiveness of our method, and achieve state-of-the-art performance<sup>1</sup>.

## 1 Introduction

Named entity recognition (NER) and relation extraction (RE) are important tasks in information extraction and knowledge base population, due to its research significance in natural language processing (NLP) and wide applications (Hosseini, 2019; Zhang et al., 2020; Qin et al., 2021; Zhang et al., 2021c). Currently, with the rapid development of multimodal learning, multimodal NER (MNER) and Multimodal RE (MRE) methods (Moon et al., 2018; Zheng et al., 2021) have been proposed to enhance linguistic representations with the aid of visual clues from images. It significantly extends the text-based models by taking images as additional inputs, since the visual contexts help to resolve ambiguous multi-sense words.

\* Corresponding Author.

<sup>1</sup>Code is available in <https://github.com/zjunlp/HVPNeT>.

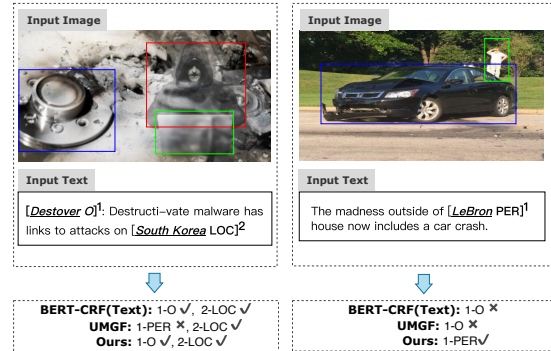


Figure 1: Motivation for robust and effective hierarchical modal fusion.

The essence of MNER and MRE tasks is how to learn great visual features and how to incorporate it into textual representation for enhancing NER and RE. Early methods (Zhang et al., 2018; Moon et al., 2018) study how to incorporate the feature of whole image into the textual representation. Yu et al. (2020); Zhang et al. (2021a); Zheng et al. (2021) further validate that object-level visual fusion is more specific and important for MNER and MRE. Recently, RpBERT (Sun et al., 2021) propose to train a classifier of whether the “Image adds to the tweet meaning” before MNER tasks. However, they heavily rely on pre-training on large extra annotated corpus of image-text relevance and only focus on the whole image with ignoring the bias of relevant object-level visual fusion. **In practice, irrelevant objects may directly exert negative effects on the text inference.** Meanwhile, it is not trivial to acquire absolutely relevant object-level visual information to enhance the text. Thus, an effective method should be derived to learn better visual representation and alleviate error sensitivity of irrelevant object images for social media NER and RE tasks.

Considering images often appear before the text in a web document, we argue that images can be regarded as the prefix for their textual descriptions,

which is inspired by prompt learning (Gao et al., 2021; Li and Liang, 2021; Liang et al., 2022; Zhang et al., 2021d) in the language model. Specifically, given a image-text pair, we prepend **object-level image feature** sequence of length  $V_i$  (visual prefix) to the text sequence at each self-attention layer of BERT (Devlin et al., 2019). Note that the visual prefix is a pluggable operation and don't require any annotation on relevance. Therefore, visual prefix can not only introduce object-level visual signals, but also further reduce the impact on the architecture representing text. Intuitively, visual prefix regarded as a prompt for text may helps alleviate the error sensitivity of irrelevant object images.

While Convolution Neural Networks (CNNs) contain the multi-scale information with pyramidal feature hierarchy (Ren et al., 2015) from low to high levels. And BERT encodes a rich hierarchy of linguistic information (Jawahar et al., 2019) from the bottom to the top. Inspired by Lin et al. (2017); Liu et al. (2018) that objects of different sizes can have appropriate feature representations at the corresponding scales, we propose to make each layer of BERT aware of hierarchical multi-scale visual features to make a more enlightened and comprehensive forecasting decision.

To this end, we propose a novel **Hierarchical Visual Prefix fusion NeTwork (HVPNeT)** for visual-enhanced entity and relation extraction. Specifically, inspired by SimVLM (Wang et al., 2021), we propose visual prefix-guided fusion mechanism involving **concatenate** object-level visual representation as the prefix of each self-attention layer in BERT, which is a more soft and robust attention module for visual enhanced NER and RE. We further design a dynamic gate for each layer to generate image-dependent paths, so that a variety of aggregated hierarchical multi-scaled visual features can be considered as visual prefix for enhancing NER and RE. Overall, we summarize the major contributions of our paper as follows:

- We present a hierarchical visual prefix fusion network towards MNER and MRE, incorporating hierarchical multi-scaled visual features through visual prefix-based attention mechanism at each self-attention layer of BERT to generate effective and robust textual representation for reducing error sensitivity.
- We utilize the exploitation of dynamic gates to fully leverage the hierarchical visual features.

Thus, textual representation of each layer in Transformer can be aware of corresponding hierarchical visual features adaptively. To the best of our knowledge, this paper is the first work to leverage hierarchical pyramidal visual features for multimodal learning.

- We evaluate our method on MNER and MRE tasks. Our experimental results on three benchmark datasets validate the effectiveness and superiority of our HVPNeT

## 2 Related work

**Multimodal Entity and Relation Extraction** As the crucial components of information extraction, named entity recognition (NER) and relation extraction (RE) have attracted much attention in the research community (Liu et al., 2019; Zhang et al., 2021b; Liu et al., 2021; Chen et al., 2021b,a). Previous studies typically focus on textual modality and standard text. As multimodal data become increasingly popular on social media platforms, early research focusing on textual modality and standard text is limited. Recently, several studies have focused on the MNER and MRE task, aiming to utilize the associate images to recognize the named entities and their relation better.

In the early stages, Zhang et al. (2018), Lu et al. (2018), (Moon et al., 2018) and Arshad et al. (2019) propose to encode the text through RNN and the whole image through CNN, then designing implicit interaction to model information between two modalities to explore multimodal NER tasks. Recently, Yu et al. (2020); Zhang et al. (2021a) propose to leverage regional image features to represent objects in the image to exploit fine-grained semantic correspondences based on Transformer and visual backbones.

While most of the current methods ignore the error sensitivity, one exception is that Sun et al. (2021), which proposes to learn a text-image relation classifier to enhance multimodal BERT to reduce the interference from irrelevant images while requiring extensive annotation for the irrelevance of image-text pairs.

### Pre-trained Multimodal Representation

The pre-trained multimodal BERT has recently achieved significant improvements in many multimodal tasks (e.g., visual question answering). We summarize and compare The existing visual-linguistic BERT models can be divided

into two aspects as follows: 1) **Architecture**. The single-stream structures consist of Unicoder-VL (Li et al., 2020), VisualBERT (Li et al., 2019), VL-BERT (Su et al., 2020), and UNITER (Chen et al., 2020b), where the text tokens and images are combined into a sequence and fed into BERT to learn contextual embeddings. The two-streams structures, LXMERT (Tan and Bansal, 2019) and ViLBERT (Lu et al., 2019), separately process the visual and language into two streams with interacting through cross-modality or co-attentional transformer layers. 2) **Pretraining tasks**. The pretraining tasks of multimodal visual-language model mainly consist of masked language modeling (MLM), masked region classification (MRC), and image-text matching (ITM). However, most of previous models are pre-trained on the datasets of image captioning (Sharma et al., 2018; Chen et al., 2015) or visual question answering where multimodal interactions are required. Applying current visual-language models to the MNER and MRE task may not result in a good performance, since **MNER and MRE mainly focus on leveraging visual information to enhance the text rather than conducting prediction on the image side.**

### 3 Methodology

As illustrated in Figure 2, we present a novel hierarchical prefix fusion network for multi-modal entity and relation extraction. Note that our method can also be applied to other visual-enhanced tasks towards text.

#### 3.1 Collection of Pyramidal Visual Feature

On the one hand, the image associated with a sentence maintains several visual objects related to the entities in the sentence, further providing more semantic knowledge to assist information extraction. On the other hand, the global image features may express abstract concepts, which play the role of a weak learning signal. **Thus, we collect multiple visual clues for multimodal entity and relation extraction, which involves taking the regional image as the vital information and the global images as the supplement.**

Given an image, we follow (Zhang et al., 2021a) to adopt the visual grounding toolkit (Yang et al., 2019) for extracting local visual objects with top  $m$  salience. Then, **we rescale the global image and object image to  $224 \times 224$  pixels as the global image  $\mathcal{I}$  and visual objects  $\mathcal{O} = \{o_1, o_2, \dots, o_m\}$ .**

In the area of CV, the feature fusion method that leveraging features from different blocks of pre-trained models (Wang et al., 2019; Kim et al., 2018; Lin et al., 2017) is widely applied for improving model performance. Inspired by such practices, we take the first step to focus on the application of pyramid features in the area of multi-modality. We propose to fuse hierarchical image features into each Transformer layer; thus, leveraging a feature pyramid is essential. Typically, given an image, we encode it with a backbone model and generate a list of **pyramidal feature maps**  $\{F_1, F_2, F_3, \dots, F_c\}$  with different scales, then map them with  $M_\theta(\cdot)$  as follows:

$$V_c = \text{Conv}_{1 \times 1}(F_c), \quad (1)$$

$$V_i = \text{Conv}_{1 \times 1}(\text{Pool}(F_i)), \quad i = 1, 2, \dots, c-1, \quad (2)$$

where  $i$  denotes the  $i$ -th block of the backbone model,  $c$  denotes the number of blocks in the visual backbone model (here is 4 for ResNet),  $\text{Pool}$  represents the pooling operation, where the features are aggregated to the same spatial sizes. **The  $1 \times 1$  convolutional layer is leveraged to map the pyramidal visual features to match the embedding size of the Transformer.**

#### 3.2 Dynamic Gated Aggregation

Although objects of different sizes can have appropriate feature representations at the corresponding scales, it is not trivial to decide which block in the visual backbone is assigned visual prefix for each layer in Transformer. To address this challenge, we propose constructing the densely connected routing space, where hierarchical multi-scaled visual features are connected with each transformer layer.

##### 3.2.1 Dynamic Gate Module

We conduct routine processes through a dynamic gate module, which can be viewed as a procedure of path decision. **The motivation of the dynamic gate is to predict a normalized vector, which represents how much to execute the visual feature of each block.** In the dynamic gate,  $g_i^{(l)} \in [0, 1]$  denotes the path probability from the  $i$ -th block of visual backbone to the  $l$ -th layer of Transformer. It is calculated as  $g^{(l)} = \mathbb{G}^{(l)}(V) \in \mathbb{R}^c$ , where  $\mathbb{G}^{(l)}(\cdot)$  denotes the gating function according to the  $l$ -th layer in Transformer,  $c$  represents the numbers of the block in backbone. We first produces the logits  $\alpha_i^{(l)}$  of the gate signals:

$$\alpha^{(l)} = f(W_l(\frac{1}{c} \sum_{i=1}^c P(V_i))), \quad (3)$$

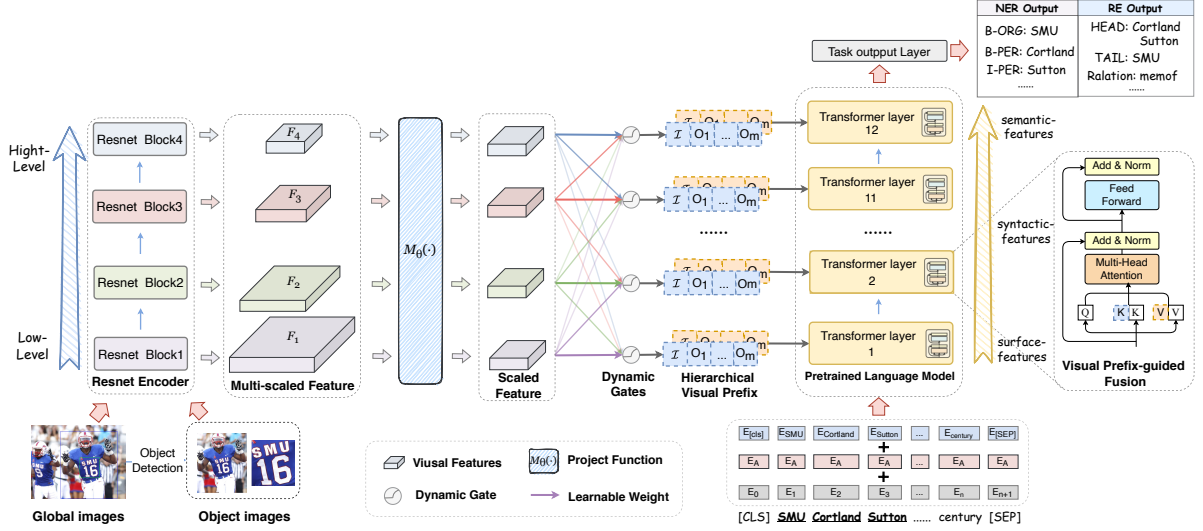


Figure 2: The overall architecture of our hierarchical visual prefix for multimodal entity and relation extraction.

where  $f(\cdot)$  denotes the activate function Leaky\_ReLU,  $P$  represents the global average pooling layer. We first squeeze the input features  $V_i$  with a shape of  $(d_i, h_i, w)$  from the  $i$ -th bloc by an average pooling operation. Then we add the features from multiple blocks to generate the average vectors. We further reduce the feature dimension by  $c$  with the MLP layer  $W_l$  and consider a soft gate via generating continuous values as path probabilities. Afterward, we generate the probability vector  $g^{(l)}$  for the  $l$ -th layer of Transformer as follows:

$$g^{(l)} = \text{Softmax}(\alpha^{(l)}) \quad (4)$$

### 3.2.2 Aggregated Hierarchical Feature

Based on the above dynamic gate  $g^{(l)}$ , we can derive the final aggregated hierarchical visual feature  $V_{gated}$  to match the  $l$ -th layer in Transformer, as:

$$V_{gated}^{(l)} = g^{(l)} V^{(l)}. \quad (5)$$

Formally, the final visual features  $\tilde{V}_{gated}^{(l)}$  corresponding to the  $l$ -th layer of Transformer is obtained by the following concatnation operation,

$$\tilde{V}_{gated}^{(l)} = [V_{gated}^{(l,1)}; V_{gated}^{(l,o_1)}; \dots; V_{gated}^{(l,o_m)}], \quad (6)$$

which will be adopted to enhance layer-level representations of textual modality through visual prefix-based attention.

### 3.3 Visual Prefix-guided Fusion

We regard hierarchical multi-scaled image feature as visual prefix, and prepend the sequence of visual

prefix to the text sequence at each self-attention layer of BERT(Devlin et al., 2019). In particular, given an input sequence  $X = \{x_1, x_2, \dots, x_n\}$ , the contextual representations  $H^{l-1} \in \mathbb{R}^{n \times d}$  is first projected into the query/key/value vector:

$$Q^l = H^{l-1} W_l^Q, K^l = H^{l-1} W_l^K, V^l = H^{l-1} W_l^V. \quad (7)$$

As for aggregated hierarchical visual features  $\tilde{V}_{gated}^{(l)}$ , we use a set of linear transformations  $W_l^\phi \in \mathbb{R}^{d \times 2 \times d}$  for  $l$ -th layer to project them into the same embedding space<sup>2</sup> of textual representation in self-attention module. Besides, we define the operation of visual prompt  $\phi_k^l, \phi_v^l \in \mathbb{R}^{hw(m+1) \times d}$  as:

$$\{\phi_k^l, \phi_v^l\} = \tilde{V}_{gated}^{(l)} W_l^\phi, \quad (8)$$

where  $hw(m+1)$  represents the length of the visual sequences,  $m$  denotes the number of visual objects detected by the object detection algorithm. Formally, the visual prefix-based attention are calculated as follows:

$$\text{Prefix\_Attention}^l = \text{softmax}\left(\frac{Q^l[\phi_k^l; K^l]^T}{\sqrt{d}}\right)[\phi_v^l; V^l]. \quad (9)$$

**Remark 1** We regard hierarchical multi-scaled visual features as visual prefix at each fusion layer and sequentially conduct multi-modal attention to update all textual states. In this way, the final textual states encode both the context and the cross-modal semantic information simultaneously. which

<sup>2</sup>Remarkably, the key and value in the self-attention module contain the different information in two types of semantic space, here 2 means that we apply two sets of transformation parameters to project aggregated visual features to match the state update process, respectively.



is beneficial to reduce error sensitivity for irrelevant object elements.

### 3.4 Classifier

Based on above description, we get the final representation of BERT,  $H^L = U(X, \tilde{V}_{gated}^{(l)})$ , where  $U(\cdot)$  denotes the operation of visual prefix-based attention. Finally, we conduct different classifier layers for NER and RE, respectively.

**Named Entity Recognition.** Following (Moon et al., 2018; Yu et al., 2020), we also adopt the CRF decoder to perform the NER task. Formally, we feed the final hidden vectors  $H^L$  of BERT to the CRF model. For a sequence of tags  $y = \{y_1, \dots, y_n\}$ , the probability of the label sequence  $y$  and the objective of NER are defined as follows (Lample et al., 2016a):

$$p(y|H^L) = \frac{\prod_{i=1}^n S_i(y_{i-1}, y_i, H^L)}{\sum_{y' \in Y} \prod_{i=1}^n S_i(y'_{i-1}, y'_i, H^L)}, \quad (10)$$

$$\mathcal{L}_{ner} = - \sum_{i=1}^M \log(p(y^{(i)} | U(X^{(i)}, \tilde{V}_{gated}))).$$

where  $Y$  represents the pre-defined label set with the BIO tagging schema, and  $S(\cdot)$  represents potential functions. Details can be referred in (Lample et al., 2016a).

**Relation Extraction.** An RE dataset can be denoted as  $\mathcal{D}_{re} = \{(X^{(i)}, r^{(i)})\}_{i=1}^M$ , the goal of RE is to predict the relation  $r \in \mathcal{Y}$  between subject entity and object entity. Specifically, a [CLS] head is utilized to compute the probability distribution over the class set  $\mathcal{Y}$  with the softmax function  $p(r|X) = \text{Softmax}(\mathbf{W}\mathbf{H}_{[CLS]}^L)$ , and the parameters of  $\mathcal{L}$  and  $\mathbf{W}$  are fine-tuned by minimizing the cross-entropy loss over  $p(r|X)$  on the entire  $\mathcal{X}$  as follows:

$$\mathcal{L}_{re} = - \sum_{i=1}^M \log(p(r^{(i)} | U(X^{(i)}, \tilde{V}_{gated}))). \quad (11)$$

## 4 Experiments

In the following section, we conduct experiments to evaluate our method on two multimodal information extraction tasks, MNER and MRE. Specifically, we adopt ResNet50 (He et al., 2016) as visual backbone and BERT-base (Devlin et al., 2019) as textual encoder. Results on three datasets demonstrate that our HVPNeT outperforms a number of unimodal and multimodal approaches.

### 4.1 Datasets

We select three datasets for our experiments: Twitter-2015 (Zhang et al., 2018) and Twitter-2017 (Lu et al., 2018) for MNER, MNRE (Zheng et al., 2021) for MRE. Statistical details of datasets and experimental details are provided in Appendix A, B.

### 4.2 Compared Baselines

We compare our HVPNeT with several baseline models for a comprehensive comparison to demonstrate the superiority of our HVPNeT. Our comparison mainly focuses on three groups of models: the text-based models, previous SOTA MNER and MRE models, and the variants of our models.

**Text-based models:** we first consider a group of representative text-based models: 1) *CNN-BiLSTM-CRF* (Ma and Hovy, 2016), 2) *HBiLSTM-CRF* (Lample et al., 2016b) and 3) *BERT-CRF* for NER. The following models are specific for RE: 4) *PCNN* (Zeng et al., 2015); 5) *MTB* (Soares et al., 2019) is an RE-oriented pretraining model based on BERT.

**Previous SOTA models:** besides, we further consider another group of previous SOTA multi-modal approaches for MNER and MRE: 1) *AdapCoAtt-BERT-CRF* (Zhang et al., 2018); 2) *OCSGA* (Wu et al., 2020); 3) *UMT* (Yu et al., 2020); 4) *UMGF* (Zhang et al., 2021a), the newest SOTA for MNER, which proposes a unified multi-modal graph fusion approach for MNER. 5) *BERT+SG* is proposed in Zheng et al. (2021) for MRE, which concatenate the textual representation from BERT with visual features generated with scene graph (SG) tool (Tang et al., 2020). 6) *MEGA* (Zheng et al., 2021), the newest SOTA for MRE, which develops a dual graph for multi-modal alignment to capture this correlation between entities and objects for better performance. 7) *VisualBERT* (Li et al., 2019), different from the above SOTA methods mainly based on co-attention, VisualBERT is a single-stream structure, which is a strong baseline for comparison. And the results of VisualBERT listed in our paper are referred from Chen et al. (2020a)

**Variants of Our Model:** we set the ablation experiments to explore the effectiveness of our design. We conduct on the same parameter settings of HVPNeT for each variant model for a fair comparison.

Modality	Methods	Twitter-2015			Twitter-2017			MNRE		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Text	CNN-BiLSTM-CRF	66.24	68.09	67.15	80.00	78.76	79.37	-	-	-
	HBiLSTM-CRF	70.32	68.05	69.17	82.69	78.16	80.37	-	-	-
	BERT-CRF	69.22	74.59	71.81	83.32	83.57	83.44	-	-	-
	PCNN	-	-	-	-	-	-	62.85	49.69	55.49
	MTB	-	-	-	-	-	-	64.46	57.81	60.86
Text+Image	AdapCoAtt-BERT-CRF	69.87	74.59	72.15	85.13	83.20	84.10	-	-	-
	OCSGA	<b>74.71</b>	71.21	72.92	-	-	-	-	-	-
	UMT	71.67	75.23	73.41	85.28	85.34	85.31	62.93	63.88	63.46
	UMGF	74.49	75.21	74.85	86.54	84.50	85.51	64.38	66.23	65.29
	BERT+SG	-	-	-	-	-	-	62.95	62.65	62.80
	MEGA	70.35	74.58	72.35	84.03	84.75	84.39	64.51	68.44	66.41
	VisualBERT	68.84	71.39	70.09	84.06	85.39	84.72	57.15	59.48	58.30
	HVPNeT-Flat	73.76	75.32	74.54	84.43	86.42	85.41	79.32	78.20	78.75
	HVPNeT-1T3	74.25	75.45	74.85	85.43	85.85	85.75	81.18	78.46	79.25
	HVPNeT-OnlyObj	74.07	76.23	75.13	85.58	87.52	86.55	81.57	80.94	81.25
	<b>HVPNeT</b>	73.87	<b>76.82</b>	<b>75.32</b>	<b>85.84</b>	<b>87.93</b>	<b>86.87</b>	<b>83.64</b>	<b>80.78</b>	<b>81.85</b>

Table 1: Performance comparison of different competitive baseline approaches for NER and RE. Since the original results of UMT, UMGF and MEGA only involve single extraction task, we reproduce their public code for more comprehensive comparison.

**HVPNeT-Flat:** This is another variant of our model without the pyramid structure. Here we assign the visual features with the output of the 4-th block of ResNet and then map the visual features to each layer corresponding to BERT to conduct image-text fusion.

**HVPNeT-1T3:** As ResNet and BERT have four blocks and 12 layers, respectively thus, it is intuitive to directly map visual features in one block to the three layers in BERT. We denote this variant as *HVPNeT-1T3* to compare with our final version with hierarchical visual features.

**HVPNeT-OnlyObj:** Visual objects are considered as fine-grained image representations. We conduct ablation by only adopting the object-level features in this model to validate the effect of the object features.

### 4.3 Overall Performance Comparison

#### 4.3.1 Main Results

The experimental results of HVPNeT and all baselines on three testing sets are presented in Table 1. From the experimental results, we can observe that:

Firstly, we can find that incorporating the visual features is generally helpful for NER and RE tasks by comparing the SOTA multimodal approaches with their corresponding text-based baselines. Despite previous multimodal approaches can generally achieve better performance, the enormous improvement of F1 score for NER is only about 2.0% (compare UMGF with BERT-CRF), which for RE is about 5.55% (compare MEGA with MTB). This observation reveals that the performance improve-

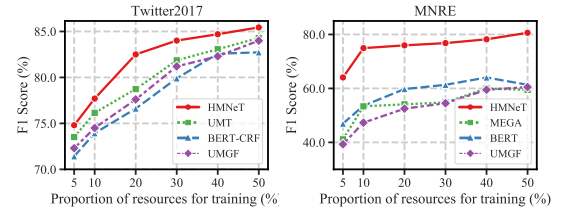


Figure 3: Performances on low-resource setting on MNER and MRE task.

ment of images on text-based NER tasks is relatively limited compared with RE tasks.

Secondly, our method is superior to the newest SOTA models UMGF and MEGA, which improves 1.36%, and 15.44% F1 scores for Twitter-2017, and MNRE datasets, respectively. It is worth noting that most of previous multimodal methods ignore the error sensitivity of irrelevant object-level images, while our method regard hierarchical visual prefix as a prompt for text. This results indicate that our method can effectively alleviate the error sensitivity irrelevant object images, which is a more robust method for visual enhanced NER and RE.

Finally, we also compare with VisualBERT, which is a pre-trained multimodal BERT with a single-stream structure. We notice that even as the pre-trained multimodal model, VisualBERT leaves much to be desired in MNER and MRE tasks, which performs worse than UMGF and MEGA, let alone our methods. We hold that VisualBERT is truly dissatisfactory since the datasets and pre-training process are less relevant to information extraction tasks.

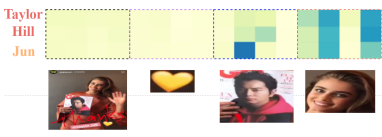


Relevant Image-text Pair	Weak Relevant Image-text Pair	Irrelevant Image-text Pair
Taylor Hill holding Jun 's GQ japan lol.	Cold front over Blyde River Canyon in Limpopo Province, South Africa.	President Bush when he sees the lights of America.
<b>Text-Images Attention of HVPNeT</b>		
		
<b>Gold Relations:</b> per/per/couple	loc/loc/contain	per/loc/place_of_residence
BERT: per / per /couple ✗	misc / misc /part_of ✗	per / loc /place_of_residence ✓
VisualBERT: per / per /peer ✓	misc / misc /part_of ✗	misc / loc /held_on ✗
MEGA: per / per /peer ✓	per / per /peer ✗	misc / loc /held_on ✗
HVPNeT(Ours): per / per /peer ✓	loc / loc /contain ✓	per / loc /place_of_residence ✓

Table 2: The first row shows the split of the relevance of image-text pairs, and the several middle rows indicate representative samples together with their entity-object attention in the test set of MNRE datasets (The y-axis represents the textual entites, and the x-axis denotes the visual objects with length of flattened 4 patches), and the bottom four rows show predicted relation of different approaches on these test samples.

### 4.3.2 Low-resource Scenario

We further conduct experiments in low-resource settings by randomly sampling 5% to 50% from the original training set to form a low-resource training set. Figure 3 shows the performance of our method in a low-resource scenario compared with several baselines. By analyzing this results, we can observe: 1) UMT and MEGA consistently outperform the compared baselines in the low-resource scenario; the improvement indicates that incorporating the visual features is still helpful for NER and RE tasks in low-resource scenarios. 2) Moreover, it can be observed that the performance of HVPNeT still outperforms the other baselines. It further proves the effectiveness and data-efficiency of our proposed method.

### 4.3.3 Cross-task Scenario

Table 3 shows performance comparison of HVPNeT and UMGF in a cross-task scenario for versatility analysis. For the first part, Twitter2017  $\rightarrow$  MNRE denotes that the trained model on Twitter-2017 is further used to train and test on MNRE. For the second part, MNRE  $\rightarrow$  Twitter-2017 represents that the trained model on Twitter-2017 is used to further train and test on Twitter-2017. From this Table, we can observe that our HVPNeT significantly outperforms UMGF by a more considerable margin. Note that our method can achieve further improvement in a cross-task scenario, while UMGF performs worse than previous results on the corresponding dataset. This justifies that our HVPNeT is robust to automatically reduce the interference of visual information of irrelevant pictures; thus,

Methods	Twitter-2017 $\rightarrow$ MNRE	MNRE $\rightarrow$ Twitter-2017
UMGF	63.85 $\rightarrow$ 62.90 $\downarrow$ (0.95)	85.51 $\rightarrow$ 84.35 $\downarrow$ (1.16)
<b>HVPNeT</b>	81.85 $\rightarrow$ 82.50 $\uparrow$ (0.75)	86.87 $\rightarrow$ 87.13 $\uparrow$ (0.26)

Table 3: Performance comparison of HVPNeT and UMGF in cross-task scenario.

more image-text data may facilitate learning better parameters for modality fusion. Besides, it is also interesting to extend our work to multi-task learning or multi-modal pre-training and we leave these for future works.

### 4.4 Detailed Model Analysis

**Ablation Study.** In this part, we conduct extensive experiments with the variants of our model to further analyze the effectiveness of our model. We illustrate the results of the variant set in Table 1. We can observe that:

(1) **Visual Prefix-guided Fusion.** The core module of our HVPNeT is visual prefix-guided fusion, which is a pluggable operation. Therefore, ablating visual prefix-guided fusion is equivalent to a purely bert-based baseline model. As shown in Table 1, HVPNeT achieve significant improvements over purely bert-based baseline model, revealing the effectiveness of pluggable visual prefix-guided fusion.

(2) **Hierarchical Visual Features.** To validate the impact of our proposed hierarchical visual features, we carry out experiments by introducing two variants: 1) HVPNeT-Flat, crudely assign single visual feature for each layer of BERT; and 2) HVPNeT-1T3, intuitively leveraging visual fea-

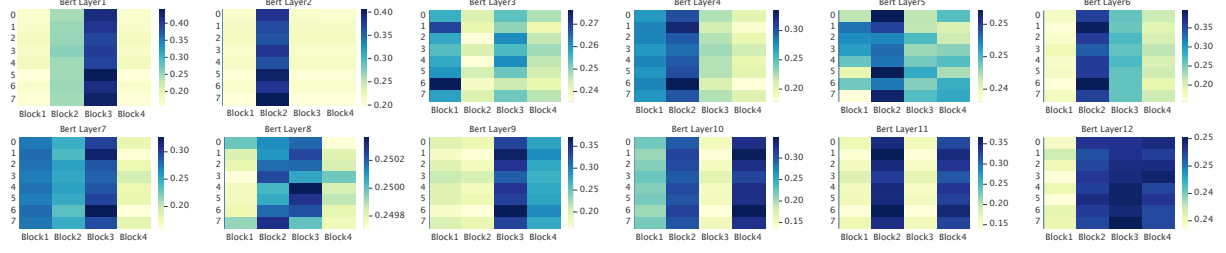


Figure 4: Visualization of dynamic gate learned on MNER task. Each subgraph denotes one layer in BERT, and the ordinate and abscissa respectively represent the instance id in a batch and the block id of ResNet.

tures from low-level to high-level blocks. We observe that HVPNeT with hierarchical visual features achieves the best performance consistently compared with the other variants. Although the HVPNeT-1T3 performs slightly lower than the version of dynamic gate, it still outperforms the crude variant HVPNeT-Flat. It reveals that the dynamic gate can automatically learn appropriate weights for multi-scaled visual representations, enabling the model to learn good visual guidance for multimodal entity and relation extraction.

(3) **Visual Clues Term.** As recent SOTA models such as UMT, UMGF, and MEGA all adopt visual objects to enhance textual representation, we conduct experiments by ablating global images to explore the impact of the visual clues. As expected, we find that HVPNeT-OnlyObj performs slightly worse than HVPNeT, which is consistent with the observation of previous works. This can be attributed to that abstract clues maybe not be associated with the text in information extraction tasks. In other words, this empirical finding demonstrates the flexibility of our methods to infuse visual clues with different granularity.

**Case Analysis for Error Sensitivity** To validate the effectiveness and robustness of our method, we conduct case analysis for image-text relevance as indicated in Table 2. We notice that VisualBERT, MEGA, and our method can recognize the relation for the relevant image-text pair. We can further find that the attention between relevant entities and objects is significant. While in the situation that image represents the abstract semantic that is weak relevant to the text, only our method success in prediction due to HVPNeT captures the more hierarchical features. It should be noted that another two multimodal baselines fail in irrelevant image-text pairs while text-based BERT and ours still predict correctly. These observations reveal that our model regards visual prefix as a prompt

for text may helps learn more robust multimodal representation, which is essential for the noise of uncorrelated object images.

**Gate Visualization** We argue that dynamic gated aggregation for hierarchical visual representation is another key component of HVPNeT achieving the superior performance. Specifically, the dynamic gated aggregation can adaptively assign different modality integration paths for different input images, thus, incorporating visual guidance with hierarchical multi-scaled information. To this end, we randomly sample eight images in a batch and visualize their gate vectors learned by HVPNeT according to 12 layers of BERT in Figure 4. Note that optimized gate vectors follow the trend of matching low-level textual semantics with low-level visual semantics and matching high-level textual semantics with high-level visual semantics. Meanwhile, the modality fusion obtained by dynamic gate learning may provide some valuable insights for efficient visual-language approaches in the future.

## 5 Conclusion and Future Work

In this paper, we propose a novel hierarchical visual prefix fusion neTwork (HVPNeT) for visual-enhanced entity and relation extraction. To be specific, we present visual prefix-guided fusion by concatenating object-level visual representation as the prefix of each self-attention layer in BERT, which is a more soft and robust attention module for visual enhanced NER and RE. We further design leveraging hierarchical multi-scaled visual representation as visual guidance for fusion. Intuitively, **Good Visual Guidance Make A Better Extractor**, and extensive experimental and results on three benchmarks have demonstrated the effectiveness and robustness of our proposed method. Meanwhile, our method also face the limitation that they don’t suitable for mulimodal tasks in visual side, such as visual grounding.



In the future, we plan to 1) explore more applications of hierarchical visual prefix in multimodal representation learning, making it more flexible and extensible; 2) try to apply the reverse version of our approach to boost visual representation with text for CV; 3) extend our approach to multitask multimodal pre-training.

## 6 Acknowledgments

We want to express gratitude to the anonymous reviewers for their hard work and kind comments. This work is funded by National Key R&D Program of China (Funding No.SQ2018YFC000004), NSFC91846204/NSFCU19B2027, Zhejiang Provincial Natural Science Foundation of China (No. LGG22F030011), Ningbo Natural Science Foundation (2021J190), and Yongjiang Talent Introduction Programme (2021A-156-G).

## References

- Omer Arshad, Ignazio Gallo, Shah Nawaz, and Alessandro Calefati. 2019. [Aiding intra-text representations with visual context for multimodal named entity recognition](#). *ArXiv preprint*, abs/1904.01356.
- Shuguang Chen, Gustavo Aguilar, Leonardo Neves, and Tamar Solorio. 2020a. [A caption is worth A thousand images: Investigating image captions for multimodal named entity recognition](#). *CoRR*, abs/2010.12712.
- Xiang Chen, Ningyu Zhang, Lei Li, Xin Xie, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021a. [Lightner: A lightweight generative framework with prompt-guided attention for low-resource NER](#). *CoRR*, abs/2109.00720.
- Xiang Chen, Ningyu Zhang, Xin Xie, Shumin Deng, Yunzhi Yao, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021b. [Knowprompt: Knowledge-aware prompt-tuning with synergistic optimization for relation extraction](#). *CoRR*, abs/2104.07650.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. 2015. [Microsoft COCO captions: Data collection and evaluation server](#). *CoRR*, abs/1504.00325.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020b. [UNITER: universal image-text representation learning](#). In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXX*, volume 12375 of *Lecture Notes in Computer Science*, pages 104–120. Springer.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 3816–3830. Association for Computational Linguistics.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society.
- Hawre Hosseini. 2019. [Implicit entity recognition, classification and linking in tweets](#). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019*, page 1448. ACM.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 3651–3657. Association for Computational Linguistics.
- Seung-Wook Kim, Hyong-Keun Kook, Jee-Young Sun, Mun-Cheon Kang, and Sung-Jea Ko. 2018. [Parallel feature pyramid network for object detection](#). In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part V*, volume 11209 of *Lecture Notes in Computer Science*, pages 239–256. Springer.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016a. [Neural architectures for named entity recognition](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 260–270, San Diego, California. Association for Computational Linguistics.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016b. [Neural architectures for named entity recognition](#). In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*,

- San Diego California, USA, June 12-17, 2016, pages 260–270. The Association for Computational Linguistics.
- Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020. [Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11336–11344. AAAI Press.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. [Visualbert: A simple and performant baseline for vision and language](#). *ArXiv preprint*, abs/1908.03557.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4582–4597. Association for Computational Linguistics.
- Xiaozhuan Liang, Ningyu Zhang, Siyuan Cheng, Zhen Bi, Zhenru Zhang, Chuanqi Tan, Songfang Huang, Fei Huang, and Huajun Chen. 2022. [Contrastive demonstration tuning for pre-trained language models](#). *CoRR*, abs/2204.04392.
- Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. 2017. [Feature pyramid networks for object detection](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 936–944. IEEE Computer Society.
- Kun Liu, Yao Fu, Chuanqi Tan, Mosha Chen, Ningyu Zhang, Songfang Huang, and Sheng Gao. 2021. [Noisy-labeled NER with confidence estimation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3437–3445. Association for Computational Linguistics.
- Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. 2018. [Path aggregation network for instance segmentation](#). In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 8759–8768. Computer Vision Foundation / IEEE Computer Society.
- Yijin Liu, Fandong Meng, Jinchao Zhang, Jinan Xu, Yufeng Chen, and Jie Zhou. 2019. [GCDT: A global context enhanced deep transition architecture for sequence labeling](#). In *Proceedings of ACL*, pages 2431–2441, Florence, Italy. Association for Computational Linguistics.
- Di Lu, Leonardo Neves, Vitor Carvalho, Ning Zhang, and Heng Ji. 2018. [Visual attention model for name tagging in multimodal social media](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1990–1999, Melbourne, Australia. Association for Computational Linguistics.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. [Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13–23.
- Xuezhe Ma and Eduard H. Hovy. 2016. [End-to-end sequence labeling via bi-directional lstm-cnns-crf](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.
- Seungwhan Moon, Leonardo Neves, and Vitor Carvalho. 2018. [Multimodal named entity recognition for short social media posts](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 852–860, New Orleans, Louisiana. Association for Computational Linguistics.
- Yujia Qin, Yankai Lin, Ryuichi Takanobu, Zhiyuan Liu, Peng Li, Heng Ji, Minlie Huang, Maosong Sun, and Jie Zhou. 2021. [ERICA: improving entity and relation understanding for pre-trained language models via contrastive learning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 3350–3363. Association for Computational Linguistics.
- Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2015. [Faster R-CNN: towards real-time object detection with region proposal networks](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. [Conceptual captions: A cleaned, hypemymed, image alt-text dataset for automatic image captioning](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 2556–2565. Association for Computational Linguistics.

- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. [Matching the blanks: Distributional similarity for relation learning](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 2895–2905. Association for Computational Linguistics.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. [VL-BERT: pre-training of generic visual-linguistic representations](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Lin Sun, Jiquan Wang, Kai Zhang, Yindu Su, and Fangsheng Weng. 2021. [Rpbert: A text-image relation propagation-based BERT model for multimodal NER](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 13860–13868. AAAI Press.
- Hao Tan and Mohit Bansal. 2019. [LXMERT: Learning cross-modality encoder representations from transformers](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China. Association for Computational Linguistics.
- Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. 2020. [Unbiased scene graph generation from biased training](#). In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 3713–3722. Computer Vision Foundation / IEEE.
- Tiancai Wang, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shabbaz Khan, Yanwei Pang, and Ling Shao. 2019. [Learning rich features at high-speed for single-shot object detection](#). In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 1971–1980. IEEE.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. [Simvlm: Simple visual language model pretraining with weak supervision](#). *CoRR*, abs/2108.10904.
- Zhiwei Wu, Changmeng Zheng, Yi Cai, Junying Chen, Ho-Fung Leung, and Qing Li 0001. 2020. [Multi-modal representation with embedded visual guiding objects for named entity recognition in social media posts](#). In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*, pages 1038–1046. ACM.
- Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. 2019. [A fast and accurate one-stage approach to visual grounding](#). In *ICCV*.
- Jianfei Yu, Jing Jiang, Li Yang, and Rui Xia. 2020. [Improving multimodal named entity recognition via entity span detection with unified multimodal transformer](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3342–3352, Online. Association for Computational Linguistics.
- Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. 2015. [Distant supervision for relation extraction via piecewise convolutional neural networks](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 1753–1762. The Association for Computational Linguistics.
- Dong Zhang, Suzhong Wei, Shoushan Li, Hanqian Wu, Qiaoming Zhu, and Guodong Zhou. 2021a. [Multi-modal graph fusion for named entity recognition with targeted visual guidance](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 14347–14355. AAAI Press.
- Ningyu Zhang, Xiang Chen, Xin Xie, Shumin Deng, Chuanqi Tan, Mosha Chen, Fei Huang, Luo Si, and Huajun Chen. 2021b. [Document-level relation extraction as semantic segmentation](#). In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 3999–4006. ijcai.org.
- Ningyu Zhang, Shumin Deng, Zhen Bi, Haiyang Yu, Jiacheng Yang, Mosha Chen, Fei Huang, Wei Zhang, and Huajun Chen. 2020. [Openue: An open toolkit of universal extraction from text](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pages 1–8. Association for Computational Linguistics.
- Ningyu Zhang, Qianghuai Jia, Shumin Deng, Xiang Chen, Hongbin Ye, Hui Chen, Huaixiao Tou, Gang Huang, Zhao Wang, Nengwei Hua, and Huajun Chen. 2021c. [Alicg: Fine-grained and evolvable conceptual graph construction for semantic search at alibaba](#). *CoRR*, abs/2106.01686.
- Ningyu Zhang, Luoqiu Li, Xiang Chen, Shumin Deng, Zhen Bi, Chuanqi Tan, Fei Huang, and Huajun Chen. 2021d. [Differentiable prompt makes pre-trained language models better few-shot learners](#). *CoRR*, abs/2108.13161.



Qi Zhang, Jinlan Fu, Xiaoyu Liu, and Xuanjing Huang. 2018. [Adaptive co-attention network for named entity recognition in tweets](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5674–5681. AAAI Press.

Changmeng Zheng, Junhao Feng, Ze Fu, Yi Cai, Qing Li, and Tao Wang. 2021. [Multimodal relation extraction with efficient graph alignment](#). In *MM '21: ACM Multimedia Conference, Virtual Event, China, October 20 - 24, 2021*, pages 5298–5306. ACM.

## A Detailed Statistics of Dataset

Dataset	Train	Dev	Test	Avg length (characters)
Twitter-2015	4,000	1,000	3,257	95
Twitter-2017	4,290	1,432	1,459	64

Table 4: Size of the datasets in numbers of tweets.

Dataset	# Sent.	# Ent.	# Rel.	# Img.
TACRED	53,791	152,527	41	-
MNRE	9,201	30,970	23	9,201

Table 5: Comparison of MNRE with existing sentence-level Relation Extraction dataset TACRED ( Sent.: sentence, Ent.: entity, Rel.: relation,Img.: image).

## B Experimental Details

This section details the training procedures and hyperparameters for each of the datasets. We use the BERT-base-uncased model from hugging face library<sup>3</sup>. We follow UMGF (Zhang et al., 2021a) to revise some wrong annotations in the Twitter-2015 dataset. Considering the instability of the few-shot learning, we run each experiment 5 times on the random seed [1, 49, 1234, 2021, 4321] and report the averaged performance. We utilize Pytorch to conduct experiments with 1 Nvidia 3090 GPUs. All optimizations are performed with the AdamW optimizer with a linear warmup of learning rate over the first 10% of gradient updates to a maximum value, then linear decay over the remainder of the training. And weight decay on all non-bias parameters is set to 0.01. We set the number of image objects  $m$  to 3. We describe the details of the training hyper-parameters in the following sections.

<sup>3</sup><https://huggingface.co/>

### B.1 Standard Supervised Setting

In the MNER task, we fix the batch size as 8 and search for the learning rates in varied intervals [1e-5, 3e-5]. We train the model for 30 epochs and do evaluation after the 16th epoch. In the MRE task, we fix the batch size as 32 and learning rates as 1e-5. We train the model for 12 epochs and do evaluation after the 8th epoch. In the two tasks, we all choices the model performing the best on the validation set and evaluate it on the test set.

### B.2 Low-Resource Setting

For different instances per class, we sample five times on the random seed [1, 2, 49, 4321, 1234] and report the averaged performance. For all models, we fix the batch size as 8 and search for the learning rates in varied intervals [3e-5, 5e-5]. We train the model for 30 epochs and do evaluation after the 16th epoch. We choose the model performing the best on the validation set and evaluate it on the test set.

### B.3 Cross-Task Setting

In the MNER task and RE task, we all use ResNet and BERT-base as the backbone, we transfer the same parameters except the classifier layer and CRF layer when we do cross-task. In further training, we fix the batch size as 8 and search for the learning rates in varied intervals [1e-5, 3e-5]. We train the model for 12 epochs and do evaluation after the 8th epoch. We choose the model performing the best on the validation set and evaluate it on the test set.