

Universal Multimodal Representation for Language Understanding

Zhuosheng Zhang[#], Kehai Chen, Rui Wang[#], Masao Utiyama, Eiichiro Sumita, Zuchao Li, Hai Zhao^{*}

Abstract—Representation learning is the foundation of natural language processing (NLP). This work presents new methods to employ visual information as assistant signals to general NLP tasks. For each sentence, we first retrieve a flexible number of images either from a light topic-image lookup table extracted over the existing sentence-image pairs or a shared cross-modal embedding space that is pre-trained on out-of-shelf text-image pairs. Then, the text and images are encoded by a Transformer encoder and convolutional neural network, respectively. The two sequences of representations are further fused by an attention layer for the interaction of the two modalities. In this study, the retrieval process is controllable and flexible. The universal visual representation overcomes the lack of large-scale bilingual sentence-image pairs. Our method can be easily applied to text-only tasks without manually annotated multimodal parallel corpora. We apply the proposed method to a wide range of natural language generation and understanding tasks, including neural machine translation, natural language inference, and semantic similarity. Experimental results show that our method is generally effective for different tasks and languages. Analysis indicates that the visual signals enrich textual representations of content words, provide fine-grained grounding information about the relationship between concepts and events, and potentially conduce to disambiguation.

Index Terms—Artificial Intelligence, Natural Language Understanding, Vision-Language Modeling, Multimodal Machine Translation.

1 INTRODUCTION

LEARNING contextualized representations of human languages is one of the major themes in natural language processing (NLP), which is also fundamental to training machines to understand human languages and handle advanced tasks, such as machine translation, question answering, and human-computer conversations. Text representation learning has evolved from standard distributed representations [1, 2] to contextualized language representation from deep pre-trained representation models (PRMs) [3, 4, 5, 6]. Despite the success of PRMs, NLP models commonly model the world knowledge (e.g. commonsense, rules, events, assertions extracted from raw texts) solely from textual features without grounding of the outside world, such as visual conception [7]. Languages are abstract

and rather difficult for the brain to retain, whereas visuals are concrete and, as such, more easily remembered [8, 9]. Adopting multimodality would be essential for better background perception. Therefore, a trend of research has been motivated to apply non-linguistic modalities to language representations [7, 10, 11, 12, 13, 14].

Most of previous works focus on joint modeling images and texts, involving vision-language (VL) pre-training [15, 16, 17, 18, 19, 20] and multimodal (MM) application tasks [14, 21, 22, 23, 24]. However, these studies rely on large-scale text-image annotations as the paired input and thus are confined to VL or MM tasks, such as image captioning and visual question answering. It is natural to boost the performance on VL and MM tasks as the concerned datasets are human-labeled with high quality. However, the essential challenge lies with the real-world scenario as there is no such high-quality annotated text-image aligned corpus for text-only NLP applications. Therefore, it is critical to investigate a general method to take advantage of visual information in a wide range of mono-modal (e.g., text-only) tasks. In addition, it is still not clear the role of images in language representation, as well as how to apply the multimodality in the standard NLP scenario.

Taking multimodal machine translation (MMT) as an example, the starting point is to leverage visual information to improve the quality of the translation from the source to the target languages. However, the effectiveness heavily relies on the availability of bilingual parallel sentence pairs with manual image annotations, which hinders the image applicability to neural machine translation (NMT). As a result, the visual information is only applied to the translation task over specific multimodal datasets [25, 26, 27, 28, 29], instead of general text-only NMT [30, 31, 32] and low-resource text-only NMT [33, 34, 35, 36]. In addition, because of the high cost of annotation, the content of one bilingual parallel

- Z. Zhang, R. Wang, Z. Li, H. Zhao are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, China and also with Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, China. K. Chen is with the School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, China. M. Utiyama and E. Sumita are with the National Institute of Information and Communications Technology (NICT), Japan. E-mail: {zhangzs, charlee}@sjtu.edu.cn; chenkehai@hit.edu.cn; {mutiyama, eiichiro.sumita}@nict.go.jp; wangrui.nlp@gmail.com; zhao-hai@cs.sjtu.edu.cn.
- H. Zhao is supported by Key Projects of National Natural Science Foundation of China (U1836222 and 61733011). R. Wang is supported by National Natural Science Foundation of China (No. 6217020129), Shanghai Pujiang Program (No. 21PJ1406800), Shanghai Municipal Science and Technology Major Project (No. 2021SHZDZX0102), Beijing Academy of Artificial Intelligence (BAAI) (No. 4), CCF-Baidu Open Fund (No. CCF-BAIDU OF2022018). K. Chen is supported by National Natural Science Foundation of China (No. 62276077) and Shenzhen College Stability Support Plan (No. GXWD20220811170358002 and GXWD20220817123150002).
- Z. Zhang and R. Wang contribute equally to this work. Part of this work was finished when Z. Zhang and Z. Li visited NICT, and R. Wang was with NICT. Corresponding author: Hai Zhao.

sentence pair is paired with a single image, which is weak in capturing the diversity of visual information. Therefore, the current study of introducing visual information falls into a bottleneck in the multimodal NMT and is not feasible for text-only NMT and low-resource NMT.

Our previous work [37] finds that using monolingual corpora with image annotations can overcome the lack of large-scale bilingual sentence-image pairs, thereby extending image applicability in NMT, with performance gains. The method of using a lookup table is task-agnostic. This work stimulates our further thinking, and we are interested in answering the three major aspects of questions:

(i) **Global Multimodality:** Can we apply the multimodality to standard text-only NLP tasks to enhance the language representations (Section 5.4), e.g., natural language generation (Section 5.1.1) and natural language understanding (Section 5.1.2)?

(ii) **Interpretability:** Why does the universal representation method work (Section 6.2)? How does multimodality improve language representation, and what is the network learned (Section 6.2-6.6)?

(iii) **Quality:** How to control the quality of visual-text alignment to reduce noise (Section 6.9)?

In this paper, we present a universal visual representation (UVR) method relying only on a seed set of task-independent annotations, instead of the existing approach that depends on large-scale task-specific image-text annotation, thus breaking the bottleneck of using visual information in standard text-only NLP tasks. For each sentence, we retrieve diverse images from either a light topic-image lookup table or pre-trained shared text-visual embedding space that is pre-trained on a large-scale of text-image pairs, to connect both the mono-modal paths of text and image embeddings. The text and images are encoded by Transformer language model (LM) and a pre-trained convolutional neural network (CNN), respectively. A simple and effective attention layer is then designed to fuse the two sequences of representations.

Our approach can be easily applied to text-only tasks without manually annotated multimodal parallel corpora. Therefore, our method is universal in terms of the task requirements, in contrast to the recent vision-language models that require large-scale and expensive annotation datasets for each downstream task. The proposed method is evaluated on 14 NLP benchmark datasets involving natural language inference (NLI), semantic similarity, text classification, and machine translation. The experiments and analysis verify the effectiveness of the proposed method. To summarize, our contributions are primarily three-fold:

(i) This work studies the universal visual representation for language representation in a broader view of the natural language processing scenario. Besides neural machine translation, this work leverages visual information as assistant signals for general NLP tasks, with the focus on investigating the global multimodality for general NLP, interpretability of effectiveness, and quality control of using universal visual representation.

(ii) For the technical side, this work proposes new methods of semantic sentence-image matching from a shared cross-modal space to give more accurately paired images as topic information. We also present a new multimodal

representation framework and systematically study the two main instances, including the model with the original TF-IDF topic-image lookup table and the newly proposed one from the retrieval from cross-modal retrieval.

(iii) Experiments are extended to 14 representative NLP tasks, which show the effectiveness of the proposed method. A series of in-depth analyses indicate that the visual signals enrich textual representations of content words, provide fine-grained grounding information about the relationship between concepts and events, and potentially conduce to disambiguation.

2 BACKGROUND

2.1 Vision-Language Integration

This study is related to that of VL methods (VLMs). Recently, there has been a great deal of interest in integrating image presentations in pre-trained Transformer architectures [15, 16, 17, 18, 19, 20, 38]. The common strategy is to take a Transformer model [32], such as BERT, as the backbone and learn aligned representations of visual and language in a pre-training manner inspired by the masked language modeling mechanism in pre-trained language models [5]. These studies require the annotation of task-dependent sentence-image pairs, which are limited to VL tasks, such as image captioning and visual question answering.

Two studies [22, 37] are closely related to general image-enhanced LM. Glyce [22] proposes incorporating glyph vectors for Chinese character representations. However, it can only be used for Chinese and only involves single image enhancement. Regarding the technical part, previous methods only benefit from one image per sentence. We propose taking advantage of a group of similar images using a filtering mechanism to form a more fine-grained visual-aware context. Our early version of this work [37] proposes using multiple images for NMT, based on a text-image lookup table trained over a sentence-image pair corpus. However, the number of images is fixed because of the lack of similarity measurement in the simple lookup method, which possibly makes the resulting model suffer from the noise of irrelevant images. This work is improved from the perspectives of motivation and technique. It is motivated by cross-modal semantic retrieval in the shared embedding space. It adopts a neural matching method with a similarity threshold to control the expected matching degree flexibly, which is generally applicable to a broader range of NLP tasks. In addition, we conduct an in-depth analysis to investigate how the visual modality helps text representation.

2.2 Visual-Semantic Embeddings

Another research line is language grounding for images whose major topic is multimodality and cross-modality between images and text. The major focus is to bridge the gap between text and images through building visual-semantic embeddings [39, 40, 41, 42, 43].

Prior studies have verified that representations of images and text can be jointly leveraged to build visual-semantic embeddings in a shared representation space [39, 40, 41, 44]. To this end, a popular approach is to connect both the mono-modal text and image encoding paths using fully connected

layers [45, 46]. The shared deep embedding can be used for cross-modal retrieval; thus, it can associate sentence text with associated images. Partly inspired by this line of research, we are motivated to incorporate visual awareness into sentence modeling by retrieving a group of images for a given sentence.

One of the first techniques to align two views of heterogeneous data is the canonical correlation analysis method [47], in which linear projections defined on both sides are optimized to maximize the cross-correlation. Recent studies have followed the two-path architecture [45, 46], in which the encoder consists of a joint embedding of textual and image representations extracted from both the images and corresponding caption. Notably, Engilberge et al. [46] adopts RNN to encode sentence embeddings in the same space with extracted image representations from CNN. Portaz et al. [48] enhances cross-modal retrieval using multilingual text. Inspired by the previous success of visual-semantic embeddings, we apply neural image retrieval from the joint space to fetch a group of associated images.

3 UNIVERSAL REPRESENTATION FRAMEWORK

This section overviews our universal representation framework. Given a sentence, we first fetch a group of matched images from our retrieval methods (details of our retrieval methods will be given in the next section). The text and images are encoded, respectively, by the text feature extractor and image feature extractor. Then the two sequences of representations are integrated using multi-head attention to form a joint representation, which is passed to downstream task-specific layers. Figure 1 overviews the whole multimodal representation model.

3.1 Encoding Layer

3.1.1 Text Encoder

We pair each sentence with the top matched m images according to the retrieval method above. Following [5], the sentence is fed into the multi-layer Transformer encoder [32] to learn the text representation $H \in \mathbb{R}^{n \times d}$ where n and d are the input text length and dimension of hidden states for the text representation.

Let $X = \{x_1, \dots, x_n\}$ be the input sentence in length n . We feed the sequence to a PRM encoder (e.g., BERT [5]). In the encoder, the input sequence is firstly mapped to embeddings. Then, the embeddings are passed to multi-head attention layers [32] to obtain the contextualized representations, which is defined as

$$H = \text{FFN}(\text{MultiHead}(K, Q, V)), \quad (1)$$

where K, Q, V are packed from the input sequence representation X . As the common practice, we set $K = Q = V$ in the implementation. MultiHead is short for multi-head attention.

3.1.2 Image Encoder

Similar to the standard way of retrieving word embeddings, the image embeddings are fetched from a lookup table \in

$\mathbb{R}^{n_m \times d_m}$ that contains the image features encoded by a pre-trained ResNet [49],¹ where n_m is the number of the total number of unique images +1 and d_m is the dimension of the image features.² The first row of the lookup table is filled by all-zero vectors, which will be used when no image is paired for the sentence.

After the feature lookup process, we obtain the image embeddings $E \in \mathbb{R}^{m \times d_m}$ for the m input images. Then, the embeddings are passed to a feedforward layer, to produce the image representation $M \in \mathbb{R}^{m \times d}$ with the same hidden dimension as H :

$$M = \text{FFN}(\text{ResNet}(E)). \quad (2)$$

There may exist cases when no word in the sentence can be found in the topic-image lookup table. When there is no paired image retrieved, we use the first-row all-zero vectors of the image lookup table as the ‘‘blank features’’ in the intuition to tell the model to ignore them.

3.2 Multimodal Integration Layer

We connect the visual and text modalities by calculating the attention between image and text features:

$$\alpha = \text{softmax}(H(W_g M + b_g)^\top), \quad (3)$$

$$H' = \alpha M, \quad (4)$$

where W_g and b_g are parameters to learn. $\alpha \in \mathbb{R}^{n \times m}$ denotes the weights assigned to the different hidden states in the sentence and the image sequences. $H' \in \mathbb{R}^{n \times d}$ is the weighted sum of all the hidden states and it represents how the sentence can be aligned to each hidden state in the image representation.

The retrieval process may possibly introduce noise of irrelevant images. To alleviate the influence, we use a neural gating mechanism for information filtering. In detail, we compute $\lambda \in [0, 1]^{n \times d}$ to weight the expected importance of image representation for each source word:

$$\lambda = \text{sigmoid}(W_\lambda H' + U_\lambda H), \quad (5)$$

where W_λ and U_λ are model parameters. We then fuse H and H' and pass the resulting representation to layer normalization and learn an effective source representation:

$$\hat{H} = \text{LayerNorm}(H + \lambda H'). \quad (6)$$

The text and image representations are jointly encoded as \hat{H} , which is fed to the task-specific layers for downstream decoding or predictions depending on task settings.

3.3 Task-specific Layer

In this section, we show how the joint representation \hat{H} is used for downstream tasks by considering NMT and NLU tasks as examples, generally following the standard procedure of the concerned tasks. For NMT, \hat{H} is directly fed to the decoder to learn a dependent-time context vector to predict the target translation. For other tasks, \hat{H} is directly fed

¹. Note that this is the standard lookup table in embedding implementations, which is not our topic-image lookup table.

². We used the maxpooling layer of ResNet, which is in the size of $\mathbb{R}^{n_m \times 2400}$.

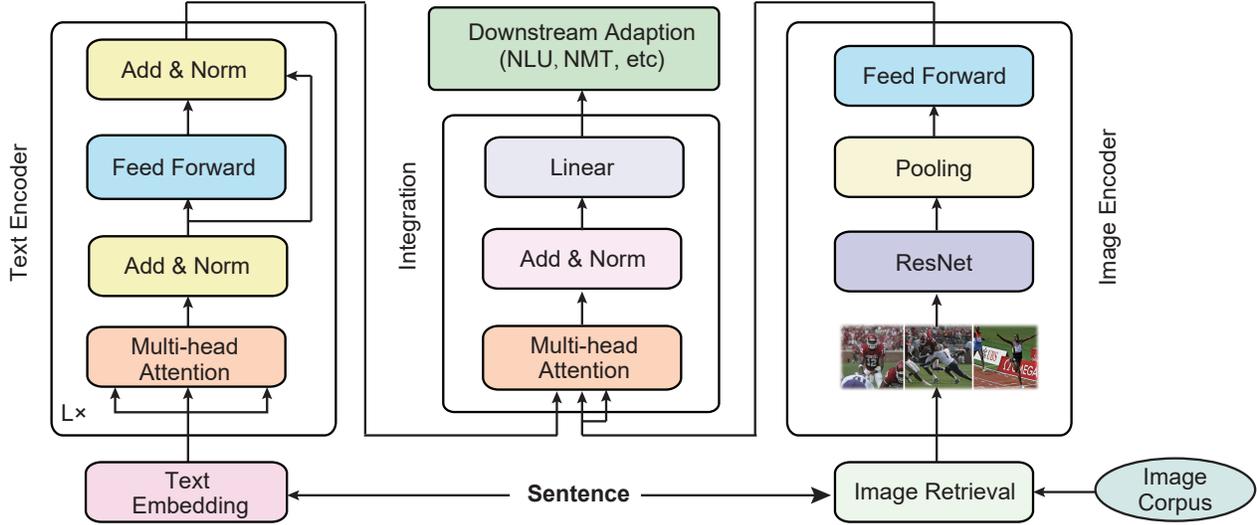


Fig. 1. Overview of the universal representation framework. Given a sentence as input, a group of related images will be retrieved by our image retrieval methods. The text and images are encoded by the text feature extractor and image feature extractor, respectively. Then the two sequences of representations are integrated using multi-head attention to form a joint representation in the same shape as the original text sequence representation. Finally, the joint representation is passed to downstream task-specific layers to give predictions.

to a feed-forward layer to make the prediction, which follows the same downstream procedure as the Transformer-based LMs, like BERT [5] and RoBERTa [50]. Specifically, for sentence-pair tasks, we maintain the pairwise input as that in LMs and separate the encoded text representation into two individual sentence representations $\{H^1, H^2\}$, according to the positions. The two text representations are integrated with the corresponding image representations respectively $\{M^1, M^2\}$, and then the resulting sequences are concatenated for prediction, $\hat{H} = \hat{H}^1 \circ \hat{H}^2$.

4 IMAGE RETRIEVAL METHODS

In this section, we describe our two visual retrieval models used for image retrieval given sentence text:

- (i) **UVR-TILT**: retrieval by topic-image lookup table;
- (ii) **UVR-CMRM**: retrieval from cross-modal embedding.

4.1 Model-I: Retrieval by Topic-Image Lookup Table

4.1.1 Topic-image Lookup Table Conversion

In this section, we will introduce the proposed universal visual representation method. Our basic intuition is to transform the existing sentence-image pairs into a topic-image lookup table,³ which assumes the topic words in a sentence should be relevant to the paired image. The procedure can be seen as the inverted index where a topic word is mapped to a list of images. Consequently, a sentence can possess a group of images by retrieving the topic-image lookup table.

To focus on the major part of the sentence and suppress the noise such as stopwords and low-frequency words, we design a filtering method to extract the “topic” words of the sentence through the term frequency-inverse document

3. We use the training set of the *Multi30K* dataset to build the topic-image lookup table.

Algorithm 1 Topic-image Lookup Table Conversion

Require: Input sentences, $S = \{X_1, X_2, \dots, X_I\}$ and paired images $E = \{e_1, e_2, \dots, e_I\}$

Ensure: Topic-image lookup table Q where each word is associated with a group of images

- 1: Obtain the TF-IDF dictionary $\mathcal{F} = \text{TF-IDF}(S)$
 - 2: Transform sentence-image pair to topic-image lookup table $Q = \text{LookUp}(S, E, \mathcal{F})$
 - 3: **procedure** TF-IDF(S)
 - 4: **for** each sentence in S **do**
 - 5: Filter stop-words in the sentence
 - 6: Calculate the TF-IDF weight for each word
 - 7: **end for**
 - 8: **return** TF-IDF dictionary \mathcal{F}
 - 9: **end procedure**
 - 10: **procedure** LOOKUP(S, E, \mathcal{F})
 - 11: **for** For each pair $\{X_i, e_i\} \in \text{zip}\{S, E\}$ **do**
 - 12: Rank and pick out the top- w “topic” words in the sentence according to the TF-IDF score in the dictionary \mathcal{F} , and each sentence is reformed as $T = \{t_1, t_2, \dots, t_w\}$
 - 13: **for** For each word t_j in T **do**
 - 14: **if** e_i not in $Q[t_j]$ **then**
 - 15: Add e_j to the corresponding image set $Q[t_j]$
 - 16: **end for**
 - 17: **end for**
 - 18: **return** Topic-image lookup table Q
 - 20: **end procedure**
-

frequency (TF-IDF),⁴ inspired by [51]. Specifically, given an original input sentence $X = \{x_1, x_2, \dots, x_I\}$ of length I and its paired image e , X is first filtered by a stopword list,⁵ and then the sentence is treated as a document g . We then

4. We describe our methods by regarding the processing unit as word though this method can also be applied to a subword-based sentence for which the subword is considered to be the processing unit.

5. <https://github.com/stopwords-iso/stopwords-en>.

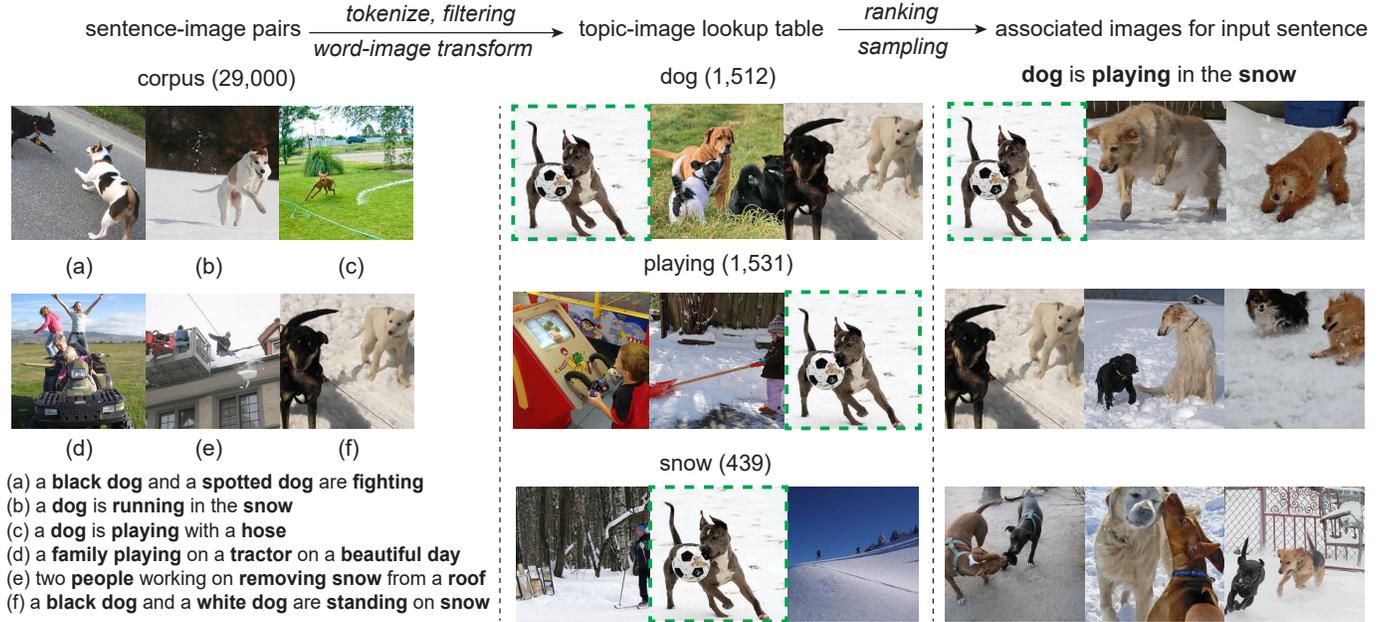


Fig. 2. Illustration of the TILT method. We first transform the existing sentence-image pairs from seed small-scale sentence-image datasets into a topic-image lookup table. For a given sentence, we extract its topic words and the associated images will be retrieved from the lookup table.

compute TF-IDF $TI_{i,j}$ for each word x_i in g ,

$$TI_{i,j} = \frac{o_{i,j}}{\sum_k o_{k,j}} \times \log \frac{|G|}{|j : x_i \in g|}, \quad (7)$$

where $o_{i,j}$ represents the number of occurrences of the word x_i in the input sentence g , $|G|$ the total number of source language sentences in the training data, and $|j : x_i \in g|$ the number of source sentences including word x_i in the training data. We then select the top- w high TF-IDF words as the new image description $T = \{t_1, t_2, \dots, t_w\}$ for the input sentence. After the preprocessing, each filtered sentence T is paired with an image e , and each word $t_i \in T$ is regarded as the topic word for image e . After processing the whole corpus (i.e., *Multi30K*), we form a topic-image lookup table \mathcal{Q} as described in Algorithm 1, in which each topic word t_i would be paired with dozens of images.

4.1.2 Image Retrieval

For the input sentence, we first obtain its topic words according to the text preprocessing method described above. Then we retrieve the associated images for each topic word from the lookup table \mathcal{Q} and group all the retrieved images together to form an image list \mathcal{G} . We observe that an image might be associated with multiple topic words so that it would occur multiple times in the list \mathcal{G} . Thus, we sort the images according to the frequency of occurrences in \mathcal{G} to maintain the same total number of images for each sentence at m .

Figure 2 illustrates the retrieval process. In the left block, we show six examples of sentence-image pairs in which the topic words are in boldface. Then we process the corpus using the topic-image transformation method demonstrated above and obtain the topic-image lookup table. For example, the word *dog* is associated with 1,512 images. For an input

source sentence, we obtain the topic words (in boldface) using the same preprocessing. Then we retrieve the corresponding images from the lookup table for each topic word. Now we have a list of images, and some images appear multiple times as they have various topics (like the boxed image in Figure 2). So we sort the retrieved image list by the count of occurrence to pick out the top- m images that cover the most topics of the sentence.

At test time, the process of getting images is done using the image lookup table built by the training set, so we do not need to use the images from the validation and test sets in *Multi30K* dataset.⁶ Intuitively, we do not strictly require the manual alignment of the word (or concept) and image, which is simpler and more general. In this way, we call our method universal visual retrieval.

4.2 Model-II: Retrieval from Cross-modal Embedding

Following Engilberge et al. [46], we train a semantic-visual embedding on a text-image corpus, which is then used for image retrieval. The semantic-visual embedding architecture comprises two paths to encode the text and images into vectors. Based on our preliminary experiments, we maintain the same settings in Engilberge et al. [46] by using the simple recurrent unit as text encoder, and the fully convolutional residual ResNet-152 [52] with Weldon pooling [53] as image encoder for our cross-modal retrieval model.

During training, each text X is paired with (i) a positive image Y that is paired with the text and (ii) a hard negative Z , which is selected as the image that has the

6. The lookup table can be easily adapted to a wide range of other NLP tasks even without any paired image, and therefore opens our proposed model to generalization.

highest similarity to the text while not being associated with it. Triplet loss [54, 55, 56] is used to enable the images to converge correctly to improve the performance of the proposed method:

$$\text{loss}(X, Y, Z) = \max(0, \gamma - E(X) \cdot E(Y) + E(X) \cdot E(Z)), \quad (8)$$

where $E(X)$, $E(Y)$, and $E(Z)$ are the embeddings of X , Y , and Z , respectively. γ is the minimum margin between the similarity of the correct caption and the unrelated caption. The loss function enables the sentence X to be closer to the corresponding image Y than the unrelated image Z . During the prediction time, the relationship between the text and images is calculated using the cosine similarity.

For general use, it is reasonable that some sentences, such as social constructs or metaphorical usage, are not paired with images after retrieval and have a low similarity score. In these cases, visual information might not be helpful. To measure how similar the retrieved images should be, we set a threshold δ to choose the top-ranked images for each sentence.

5 EXPERIMENTS

5.1 Task Settings

Our evaluation is performed on the widely-used natural language generation and understanding tasks involving 14 NLP benchmark datasets that involve machine translation, natural language inference (NLI), semantic similarity, and text classification. Part of the NLU tasks is available from the GLUE benchmark [57], which is a collection of nine NLU tasks.

5.1.1 Neural Machine Translation

Five widely-used translation tasks are used for model evaluation, including WMT’16 English-to-Romanian (En-Ro), WMT’14 English-to-German (En-De), WMT’14 English-to-French (En-Fr), and Multi30K dataset for WMT’16 and WMT’17, which are standard corpora for NMT and MMT evaluation.

(i) For the En-Ro task, we experiment with the officially provided parallel corpus: Europarl v7 and SETIMES2 from WMT’16 with 0.6M sentence pairs. We use *newsdev2016* as the validation set and *newstest2016* as the test set.

(ii) The En-De task has 4.43M bilingual sentence pairs of the WMT14 dataset used as training data, including Common Crawl, News Commentary, and Europarl v7. The *newstest2013* and *newstest2014* datasets are used as the validation set and test set, respectively.

(iii) The En-Fr task has 36M bilingual sentence pairs from the WMT14 dataset used as training data. *Newstest12* and *newstest13* are combined for validation and *newstest14* is used as the test set, following the setting of [31].

(iv) *Multi30K* dataset contains 29K English \rightarrow {German, French} parallel sentence pairs with visual annotations. The 1,014 English \rightarrow {German, French} sentence pairs with visual annotations serve as the validation set. For WMT’16 and WMT’17 tasks, we have two test sets, *test2016* and *test2017*, with 1,000 pairs for each.

5.1.2 Natural Language Understanding

The NLU task involves natural language inference, semantic similarity, and classification subtasks.

Natural Language Inference involves reading a pair of sentences and assessing the relationship between their meanings, such as entailment, neutral, and contradiction. We evaluate the proposed method on four diverse datasets: SNLI [58], MNLI [59], QNLI [60], and RTE [61].

Semantic Similarity aims to predict whether two sentences are semantically equivalent. Three datasets are used: Microsoft Paraphrase Corpus (MRPC) [62], Quora Question Pairs (QQP) dataset [63], and Semantic Textual Similarity benchmark (STS-B) [64].

Classification CoLA [65] is used to predict whether an English sentence is linguistically acceptable. SST-2 [66] provides a dataset for sentiment classification that needs to determine whether the sentiment of a sentence extracted from movie reviews is positive or negative.

5.2 Retrieval Setup

This part describes the implementation of the image retrieval by the topic-image lookup table (TILT) and cross-modal retrieval model (CMRM):

TILT: We segment the sentences using the same BPE vocabulary as that for each source language. We select top-8 ($w = 8$) high TF-IDF words, and the default number of images m is set to 5.⁷ The detailed case study is shown in Section 6.9. Image features are extracted from the averaged pooled features of a pre-trained ResNet50 CNN [49]. The dimension of the feature maps is $V \in R^{2048}$.

CMRM: The cross-modal retrieval model is trained on the MS-COCO dataset [67], which contains 123,287 images with five English captions per image. It is split into 82,783 training images, 5,000 validation images, and 5,000 test images. We use the Karpathy split [40] that forms 113,287 training, 5,000 validation and 5,000 test images. The model is implemented following the same settings as Engilberge et al. [46], and produces state-of-the-art results (94.0% R@10) for cross-modal retrieval. To ensure that each task can enjoy enough images, we set the similarity threshold δ to 0.4 and rank the paired images according to the similarity score. The maximum number of retrieved images m for each sentence is set to eight according to our preliminary experiments.

Multi30K and COCO datasets are used as the candidate seed image retrieval corpus for our downstream tasks.

5.3 Model Implementation

Since our task involves text generation and understanding, we have two kinds of baselines, the NLG model for translation and the NLU model for the other tasks.

5.3.1 NLG Model

Our baseline for NLG is encoder-decoder Transformer [32]. We use six layers for the encoder and the decoder. The number of dimensions of all input and output layers is set to 512 and 1024 for *base* and *big* models. For MMT experiments

⁷ In some cases when there is no paired image retrieved, we use the first-row all-zero vectors of the image lookup table as the “blank features”.

TABLE 1

Results for the NMT tasks. “++/+” after the BLEU score indicates that the proposed method (base: 5-8; large:9-12) was significantly better than the corresponding baseline Transformer (base or big) at significance level $p < 0.01/0.05$.

#	Model	En→Ro		En→De		En→Fr	
		BLEU	#Param	BLEU	#Param	BLEU	#Param
<i>Text-only Transformer</i>							
1	Transformer-Base	32.66	61.54M	27.31	63.44M	38.52	63.83M
2	Transformer-Big	33.85	207.02M	28.45	210.88M	41.10	211.66M
<i>Our MMT systems</i>							
3	UVR-TILT _{Multi30K}	33.78++	63.04M	28.14++	64.94M	39.64++	65.33M
4	UVR-TILT _{COCO}	34.08++	63.04M	27.79+	64.94M	39.84++	65.33M
5	UVR-CMRM _{Multi30K}	34.38++	63.04M	27.82+	64.94M	39.76++	65.33M
6	UVR-CMRM _{COCO}	34.40++	63.04M	27.86+	64.94M	40.24++	65.33M
7	UVR-TILT _{Multi30K}	34.46+	211.02M	29.14++	214.89M	41.83+	215.66M
8	UVR-TILT _{COCO}	34.51+	211.02M	29.18++	214.89M	41.76+	215.66M
9	UVR-CMRM _{Multi30K}	34.60++	211.02M	28.96+	64.94M	41.79+	215.66M
10	UVR-CMRM _{COCO}	34.62++	211.02M	29.21++	64.94M	41.82+	215.66M

TABLE 2

Results (BLEU) from the test2016 and test2017 for the MMT task. “++/+” after the BLEU score indicates that the proposed method (base: 5-8; large: 9-12) was significantly better than the corresponding baseline Transformer (base or big) at significance level $p < 0.01/0.05$.

#	Model	En-De			En-Fr		
		Test2016	Test2017	#Param	Test2016	Test2017	#Param
<i>Text-only Transformer</i>							
1	Transformer-Base	35.59	26.31	49.15M	57.88	48.55	49.07M
2	Transformer-Big	36.86	27.62	186.38M	56.97	48.17	186.23M
<i>Standard MMT systems</i>							
3	MMT-Base	35.09	27.10	50.72M	57.40	48.02	50.65M
4	MMT-Big	35.60	28.02	190.58M	57.87	49.63	190.43M
<i>Our MMT systems</i>							
5	UVR-TILT _{Multi30K}	35.72	26.87+	50.72M	58.32+	48.69	50.65M
6	UVR-TILT _{COCO}	35.67	26.89+	50.72M	58.21+	48.73	50.65M
7	UVR-CMRM _{Multi30K}	36.38+	27.34++	50.72M	58.53+	49.28+	50.65M
8	UVR-CMRM _{COCO}	35.78	26.92+	50.72M	58.46+	48.58	50.65M
9	UVR-TILT _{Multi30K}	37.02	28.63++	190.58M	57.53+	48.46	190.43M
10	UVR-TILT _{COCO}	36.94	28.69++	190.58M	57.62+	48.39	190.43M
11	UVR-CMRM _{Multi30K}	37.16	28.82++	190.58M	58.37++	48.77+	190.43M
12	UVR-CMRM _{COCO}	37.28+	28.71++	190.58M	57.60+	48.42	190.43M

on the Multi30K dataset, we also use the tiny setting where the dimension of the input and output layer is 128. The inner feed-forward neural network layer is set to 2048. The heads of all multi-head modules are set to eight in both the encoder and decoder layers. For the *Multi30K* dataset, we further evaluate a multimodal baseline (denoted as MMT) where each source sentence was paired with an original image. The other settings were the same as our proposed model.

The byte pair encoding algorithm is adopted to segment sentences into subword sequences, with the vocabulary size set to 40,000. In each training batch, a set of sentence pairs contains approximately 4096×4 source tokens and 4096×4 target tokens. During training, the value of label smoothing is set to 0.1, and the attention dropout and residual dropout rates are $p = 0.1$. We used Adam optimizer [74] to tune the parameters of the model. The learning rate is varied under a warm-up strategy with 8,000 steps. For evaluation, we validate the model with an interval of 1,000 batches on the validation set. For the *Multi30K* dataset, we train the model

up to 10,000 steps, and the training will be early-stopped if the validation set BLEU score does not improve for ten epochs. For the En-De, En-Ro, and En-Fr tasks, following the training of 200,000 batches, the model with the highest BLEU score of the validation set is selected to evaluate the test sets. During the decoding, the beam size is set to five. Multi-bleu.perl is used to compute case-sensitive 4-gram BLEU scores for all test sets.⁸ We follow the model configurations of [32] to train big models for WMT En-Ro, En-De, and En-Fr translation tasks. The experiments of NLG are conducted with *fairseq* [75].⁹

For the statistical tests, we perform the paired bootstrap resampling test [76] to measure the reliability of the conclusion that our system is better than the baseline. Our im-

8. <https://github.com/moses-smt/mosesdecoder/tree/RELEASE-4.0/scripts/generic/multi-bleu.perl>.

9. <https://github.com/pytorch/fairseq>.

TABLE 3

Comparison with public methods on the Multi30K MMT dataset. The results of existing methods are from [68]. “++/+” after the BLEU score indicates that the proposed method was significantly better than the corresponding baseline Transformer (tiny) at significance level $p < 0.01/0.05$.

#	Model	En-De			En-Fr		
		Test2016	Test2017	#Param	Test2016	Test2017	#Param
<i>Text-only Transformer</i>							
1	Transformer-Tiny	40.38	32.86	2.6M	61.00	52.42	2.6M
<i>Existing MMT systems</i>							
2	GMNMT [69]	39.8	32.2	4.0M	60.9	53.9	-
3	DCCN [70]	39.7	31.0	17.1M	61.2	54.3	16.9M
<i>Our MMT systems</i>							
4	UVR-TILT _{Tiny}	41.27++	33.62++	2.9M	61.60+	54.83++	2.9M
5	UVR-CMRM _{Tiny}	40.94+	33.11+	2.9M	61.50+	53.64++	2.9M

TABLE 4

Test results on the GLUE benchmark. The best results are marked in boldface.

#	Model	Classification		Semantic Similarity			Language Inference				Average
		CoLA	SST-2	MRPC	STS-B	QQP	MNLI	QNLI	RTE	SNLI	
<i>Public Systems</i>											
1	BERT [5]	60.5	94.9	85.4	87.6	89.3	86.7	92.7	70.1	-	83.4
2	MT-DNN [71]	62.5	95.6	88.2	89.5	89.6	86.7	93.1	81.4	91.6	86.4
3	BERT + Voken-cls [72]	-	92.2	-	-	88.6	82.6	88.6	-	-	-
4	UniT [73]	-	91.5	-	-	88.4	79.8	88.0	-	-	-
<i>Our Systems</i>											
5	Baseline (BERT _{WWM})	63.6	93.6	87.0	90.2	88.8	87.2	93.9	77.3	91.6	85.9
6	UVR-TILT _{Multi30K}	62.5	94.7	87.7	89.8	89.4	87.2	94.1	84.5	91.7	86.8
7	UVR-TILT _{CoCo}	62.8	94.9	87.4	90.2	89.7	86.9	94.0	83.6	91.7	86.8
8	UVR-CMRM _{Multi30K}	63.0	94.3	87.8	90.2	89.6	87.3	93.8	83.8	91.6	86.8
9	UVR-CMRM _{CoCo}	63.2	94.6	87.9	90.3	89.8	87.4	94.2	83.9	91.7	87.0

plementation is based on the public toolkit.¹⁰ Two thousand bootstrap samples are used for each significance test.

5.3.2 NLU Model

For the NLU tasks, the baseline is encoder-only BERT [5].¹¹ We use the whole-word-mask (WWM) version of the pre-trained weights due to its more stable, reproducible, and slightly better performance than the original large version [5]. The initial learning rate is set in the range $\{2e-5, 3e-5\}$ with a warm-up rate of 0.1 and L2 weight decay of 0.01. The batch size is selected from $\{16, 24, 32\}$. The maximum number of epochs is set in the range $[2, 5]$. Texts are tokenized using SentencePiece,¹² with a maximum length of 128.

5.4 Main Results

Tables 1-4 show the results for the 14 NMT, MMT, and NLU tasks, respectively. According to the results, we have the following observations:

(i) According to the machine translation results in Tables 1-2, the proposed UVR methods significantly outperform the baselines according to the statistical test, demonstrating

the effectiveness of modeling visual information for text-only NMT. In particular, the superiority is observed in the translation tasks of three language pairs with different training data scales, verifying that the proposed approach is a universal method for improving translation performance.

(ii) Our method introduces only 1.5M and 4.0M parameters for the base and big Transformers, respectively. The number is less than 3% of the baseline parameters as we use the fixed image embeddings from the pre-trained ResNet feature extractor. Besides, the training time is basically the same as the baseline model (Section 6.10).

(iii) Results in Tables 2-3 show that our model can generally outperform the Transformer baseline in multi-modal settings that could benefit from the gold sentence-image annotations. Compared with the results in text-only NMT, we find that the image enhancement sometimes gives marginal contribution, which is consistent with the findings in previous work [77, 78, 79]. The most plausible reason might be that the sentences in *Multi30K* are quite simple, short, and repetitive, so that the source text itself is sufficient to perform the translation [10, 79]. We also see that the big models sometimes show inferior results. The possible reason is that the dataset is too small to effectively train such big models, which easily suffer from over-fitting issues. The hypothesis is also supported by our superior results with the tiny model setting in Table 3. The observation verifies our assumption of the current bottleneck of MMT due to

10. <https://github.com/neubig/util-scripts/blob/master/paired-bootstrap.py>

11. <https://github.com/huggingface/transformers>.

12. <https://github.com/google/sentencepiece>.

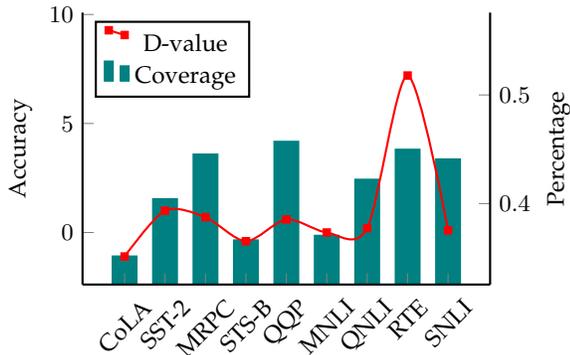


Fig. 3. Accuracy difference between our method and baseline compared with the coverage percentage of tokens that can be paired with images in each dataset.

TABLE 5

Validation results on GLUE datasets in different sizes: small datasets with less than $10k$ examples (RTE and STS-B), and a large dataset with more than $10k$ examples (QNLI). The MT-DNN results are reproduced using the released weights [71].

Model	RTE	STS-B	QNLI
MT-DNN _{base}	78.94±0.83	88.14±0.40	90.32±0.12
w/ UVR-TILT	81.59±0.63	90.16±0.07	91.31±0.19
MT-DNN _{large}	78.70±1.65	90.16±0.17	92.04±0.27
w/ UVR-TILT	82.19±1.37	91.32±0.12	92.78±0.12

the limitation of *Multi30K* and shows the necessity of our new methodology of transferring multimodality into more standard and mature text-only NMT tasks.

(iv) Table 4 shows our method is generally helpful for a wide range of NLU tasks in the GLUE benchmark, which verifies the effectiveness of modeling visual information for language understanding. We are interested in whether public methods, such as MT-DNN, can be further enhanced by our method, we apply our UVR-TILT method to the MT-DNN model based on the same implementation in BERT. According to the results in Table 5, both the base and large models are enhanced, and we observe consistent gains in different datasets, especially the small datasets. For the results in Tables 4-5, we notice a few inferior or marginally better performances in the CoLA, MNLI, and QNLI tasks. We calculate the accuracy difference (D-value) between our method and baseline compared with the coverage percentage of tokens that can be paired with images in each dataset using UVR-TILT_{Multi30K}. From Figure 3, we see that those datasets are commonly paired with a relatively small number of images so that the visual signals can enhance only a small proportion of token representations. In addition, some datasets, i.e., CoLA, mainly require linguistic knowledge for solving the tasks, so introducing visual modality might not benefit such tasks, which corresponds to the common shortcoming of vision injection for language tasks. For MNLI and QNLI, the possible reason for the marginal improvements would be that both of the datasets are quite large. Still, the task is relatively simple, so the model might well solve the tasks directly via the text representations. In this scenario, the visual features might only provide the regularization effect to improve the model robustness [80, 81, 82].

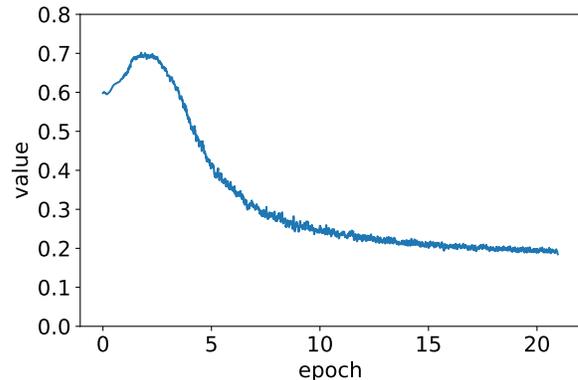


Fig. 4. Illustration of the gate values λ with the UVR-TILT method on Multi30K En-De Test2016.

(v) For the two retrieval methods, we observe that *UVR-CMRM* is slightly better than *UVR-TILT* in general, and the results of the two methods are pretty close for the MMT task. We find a performance tradeoff between them: there might be more accurate similarity calculation after cross-modal pre-training than the direct topic extraction. However, controlling the proper similarity threshold would be a heuristic for each dataset. In contrast, the advantage of *UVR-TILT* is the simple preprocessing, which only requires TF-IDF-based topic extraction and matching.

(vi) We also compare the results using different seed corpora, i.e., Multi30K and COCO. For the MMT task in Table 2, using Multi30k is basically similar to or slightly better than COCO in general, as the task could enjoy the images in the same domain. For the out-of-domain evaluations in Table 1 and Table 4, we observe that using COCO generally achieves better results because the size of COCO images is three times that of Multi30K, which could provide more diverse image features.

6 ANALYSIS

This section presents our exploration on the role of visual contexts, which involves two aspects, when the visual context helps and how the visual context helps language representations. In the following analysis part, we use Multi30K for UVR-TILT and COCO for UVR-CMRM by default.

6.1 Dynamics of the visual information

To explore the role of visual context in the training process, we illustrate the gate values λ (defined in Eq. 5) in Figure 4 where a larger value indicates more dependence on the visual context in the fusion process. We observe that the model relies on the visual context in the early stages, and the role of visual information changes dynamically across training. When the training starts, the model accommodates the visual information to a large extent ($\lambda \geq 0.6$), indicating that the model tends to trust the visual context, which could provide useful information in the early stages. With more knowledge captured as the training continues, the contributions of the visual contexts appeal to decay. In the following parts, we will further discuss the specific effectiveness of visual representations in language modeling.



A girl in a purple tutu dances in the yard.

A little girl is walking over a path of numbers.

A girl jumping rope on a sidewalk near a parking garage.

A young girl washes an automobile.

Fig. 5. Examples of sentences that share the same retrieved images, in which the common topic is about “girl”. Sentences with similar topics tend to be paired with similar or even the same images, and vice versa. This means that images may provide topic information, which benefits the modeling of similar sentences.

6.2 Pairwise relationship across modalities

The benefits of the universal representation method could be two folds: (i) the content connection of the sentences and images; (ii) the topic-aware co-occurrence of similar images and sentences. According to Distributional Hypothesis [83], which states that *words that occur in similar contexts tend to have similar meanings*, we are inspired to extend the concept in the multimodal world, *the sentences with similar meanings would be likely to pair with similar even the same images*. Therefore, the consistent images (with a related topic) could play the role of topic or type hints for similar sentence modeling.

After using our image retrieval method, sentences with similar topics tend to be paired with similar or even the same images, and vice versa. Figure 5 shows examples in which the common topic is about “girl”. This means that images may provide topic information, which benefits the modeling of similar sentences. Thus, aside from the image embeddings’ inner meaning (vectors), there is a mapping relationship between the sentence and images after the retrieval process.

For the image embeddings, as described in Section 3.1.2, we adopt the embedding lookup to fetch the embedding features for each image, which is very similar to the way of using word embedding by treating each image as a “word”. The weights of the embedding features are derived from the average pooled output of ResNet, where each image is represented as a 2400-d vector. For all the 29,000 images (e.g., using Multi30K), we have an embedding layer with size (29000, 2400). The “content” of the image can be seen as embedding initialization. The pre-initialized embedding weights might yield slight improvement gains. However, the neural network can also be effectively trained with random initialization [84, 85]. In contrast, whether to use the embedded feature is more critical. In other words, the mapping relationship of the sentences and images in image embedding would be essential, i.e., similar sentences (with the same topic words) tend to map the same or similar image after the word-image lookup process.

To verify the hypotheses, we conduct the following ablations: we replace the ResNet50 feature extractor in our UVR-TILT model with (1) *ResNet101* and (2) *ResNet152*; additionally, we compare the results with the following operations: (3) *Shuffle*: shuffle the image features but retain

TABLE 6

Ablation for the image embedding operation on En-Ro. The scores are reported by means and standard deviations for three random seeds.

Method	BLEU Score
Baseline	32.75±0.10
UVR-TILT	33.72±0.08
w / (1) Res101	33.65±0.06
w / (2) Res152	33.82±0.07
w / (3) Shuffle	33.40±0.14
w / (4) Random Init	33.08±0.17
w / (5) Random Mapping	32.05±0.14

the lookup table; (4) *Random Init*: randomly initialize the image embedding but keep the lookup table; (5) *Random Mapping*: randomly retrieve unrelated images.

Table 6 shows the ablation results. The BLEU scores of models 1-4 are close to the proposed UVR method, and those ablated methods still outperform the baseline, indicating that using image features generally yields better performance than the baseline. In addition, either replacing the trained image features (model 4) or disturbing the mapping information (model 5) leads to a performance drop ($\downarrow 0.64/\downarrow 1.67$, respectively), which indicates that both the image features and mapping information are contributing factors. Compared with image features, the mapping information has a larger impact, which verifies our prior hypothesis that the consistent images with a related topic could play the role of topic or type hints for similar sentence modeling in the whole training process. With the mapping information, the same images will be assigned to the same context. During training, the image features will be learned just like word embeddings. Therefore, it does not mean that the image features are not very helpful, but the mapping information in the lookup table reduces the dependence on the trained image features.

From the view of an individual image, image content (embedding) has an effect. If we maintain the pairwise relationship between the sentence and image, the result is still higher than the baseline, even with shuffled or random image embeddings. This indicates that the pairwise relationship is a vital contributor. From the macro perspective of sentence-image co-occurrence, image information plays the role of topic information, where similar sentences tend

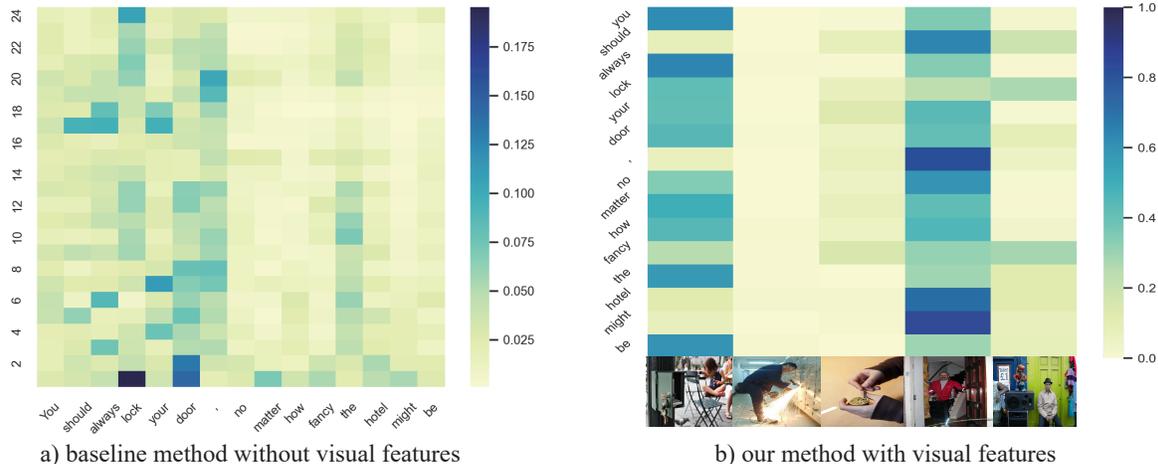


Fig. 6. Visualization of (a) attention weights of the input tokens with regards to the probed token “lock” across different layers (Y-axis) using the BERT baseline; (b) image-to-word attention from our model. The illustration shows that the images provide fine-grained grounding information about the relationship between concepts and events, e.g., “lock”, “door”, “fancy”, “hotel”.

TABLE 7

Results of MMT with incomplete source texts by removing visually grounded tokens in the Multi30K dataset. The scores are reported by means and standard deviations for three random seeds.

Model	En-De		En-Fr	
	Test2016	Test2017	Test2016	Test2017
Baseline	10.94±0.21	7.75±0.24	18.61±0.16	15.01±0.15
UVR-TILT	12.80±0.18	9.15±0.19	19.60±0.17	15.77±0.20

to pair with similar images. The observation corresponds to the distributional hypothesis. This explains the potential effects of the pairwise relationship.

This finding may potentially facilitate future research because most existing studies focus on the content of the individual image itself. We highlight the pairwise relationship across modalities as a different research line to bridge the gap between language and image modeling.

6.3 Handling incomplete source texts

Content words are naturally related to specific content, such as car, room, play. We collect a list of tokens that have more than ten occurrences in the Multi30K training set after removing all stop words following [72]. We remove those tokens in the source sentence, which occupy 42.87% of tokens in the token dictionary. Table 7 shows the results of the baseline model and our model with the incomplete source texts. We observe that our model achieves noticeable gains over the baseline. The results verify that the visual representation can reduce the gap of the missing information from the content words in the source texts.

6.4 Knowledge grounding with the visual context

To gain an insight into the process of multimodal integration by our model, we analyze the attention distributions (α in Eq.3) at the multimodal integration layer. Figure 6 shows the attention distributions of (a) the baseline and (b) our model¹³ for an example randomly selected from

13. Our model is the UVR-TILT trained on the CoLA dataset.

TABLE 8

Results (BLEU score) of the multimodal disambiguation experiments on WAT’19 English to Hindi dataset. The scores are reported by means and standard deviations for three random seeds.

Model	Validation	Test	Challenge
Baseline	47.04±0.27	39.33±0.24	20.52±0.14
UVR-CMRM	47.49±0.06	39.81±0.12	21.62±0.08

our GLUE validation sets, “You should always lock your door, no matter how fancy the hotel might be.” For comparison, the baseline is implemented following Abnar and Zuidema [86]. Concretely, we collect the attention weights of all the input tokens with regards to a targeted token “lock” across different layers (i.e., $\{2, 4, \dots, 24\}$). As the baseline uses the representation of the last layer for prediction, we focus on the attention distributions of the last layer.

Compared with the baseline that only captures partial relations in the last layer, e.g., with a lack of relationship among {“lock”, “door”, “hotel”}, our model provides more fine-grained connections. In detail, two generic patterns are observed in these examples.

(i) the images appear to match the concepts and actions with the texts, in other words, the images tend to provide fine-grained grounding information about the relationship between concepts and events, e.g., {“lock”, “door”, “fancy”, “hotel”}.

(ii) our model can resist irrelevant information from noisy images. For example, the second and third images yield low attention scores for the texts.

6.5 Disambiguation

A natural intuition of using visual clues for text representation is the advantage of alleviating the ambiguity of language. To evaluate the model performance for disambiguation, we use a dataset from the HVG [88], which serves as a part of the WAT’19 Multimodal Translation Task.¹⁴ The dataset consists of a total of 31525 randomly selected images

14. <http://lotus.kuee.kyoto-u.ac.jp/WAT/WAT2019/index.html>

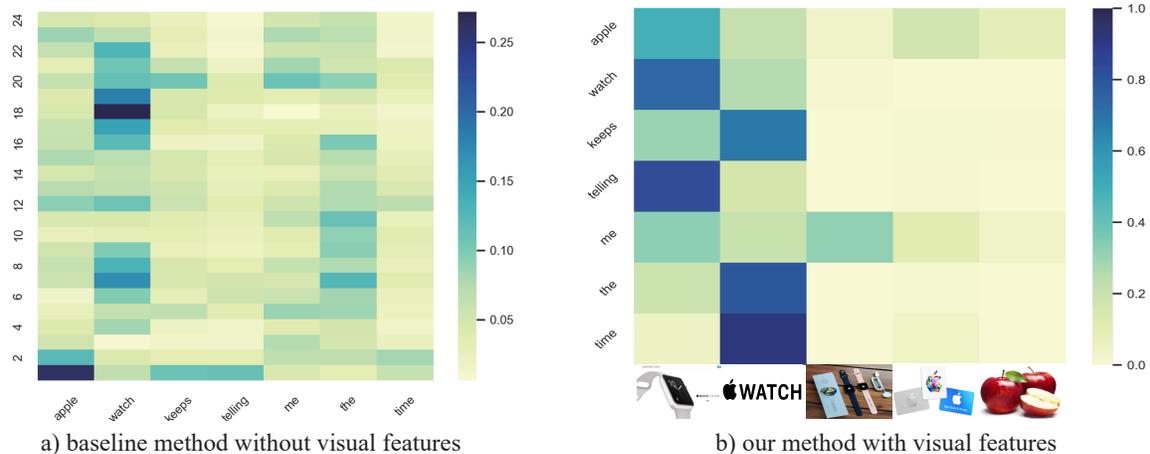


Fig. 7. Visualization of (a) attention weights of the input tokens with regards to the ambiguous token “apple” across different layers (Y-axis) using the BERT baseline; b) image-to-word attention from our model. The illustration shows that the images bridge the connection between “apple”, “watch”, “time”, helping disambiguate the meaning of “apple”.

TABLE 9
Selected eight probing tasks [87] to study what syntactic and semantic properties are captured by the encoders.

Probing Tasks		Content
Syntactic	TrDep	Checking whether an encoder infers the hierarchical structure of sentence
	ToCo	Sentences should be classified in terms of the sequence of top constituents immediately below the sentence node
	BShif	Testing whether two consecutive tokens within the sentence have been inverted
Semantic	Tense	Asking for the tense of the main clause verb
	SubN	Focusing on the number of the main clause ’ s subject
	ObjN	Testing for the number of the direct object of the main clause
	SoMo	Some sentences are modified by replacing a random noun or verb with another one and the classifier should tell whether a sentence has been modified
	CoIn	Containing sentences made of two coordinate clauses

TABLE 10
Classification accuracy on eight probing tasks of evaluating linguistics embedded in the encoder outputs.

Model	Syntactic			Semantic				
	TrDep	ToCo	BShif	Tense	SubN	ObjN	SoMo	CoIn
Baseline	28.34	58.33	76.34	80.66	72.02	68.57	64.42	67.51
UVR-CMRM	28.53	58.64	77.72	80.97	73.79	69.66	65.44	67.23

from Visual Genome [89] and a parallel image caption corpus in English-Hindi for selected image segments. The training part consists of 29K English and Hindi short captions of rectangular areas in photos of various scenes, and it is complemented by three evaluation subsets: validation, test, and challenge test set (Challenge). The challenge test set is created by searching for (particularly) ambiguous English words based on the embedding similarity and manually selecting those where the image helps to resolve the ambiguity. We do not use the images but follow the same settings as the experiments on Multi30K. As the results shown in Table 8, we observe that our UVR model works effectively on the challenge disambiguation set, indicating that the visual information induced by retrieved images allows disambiguation of translation.

Figure 7 shows a heatmap of the attention visualization on an ambiguous sentence, “apple watch keeps telling me the time”. In Figure 7(a), the baseline model fails to capture the relationship between “apple” with “time”. In Figure

7(b), the retrieved images bridge the connection among “apple”, “watch” and “time”, which helps disambiguate the meaning of “apple”.

6.6 Linguistic Analysis

We are interested in what knowledge is learned in the universal representations. In this section, we select eight widely-used language probing tasks [87] (see Table 9) to study what kind of syntactic and semantic properties are captured by the encoders. Specifically, we use the encoders of the baseline BERT-based SNLI model,¹⁵ and our UVR-CMRM visual representation model to generate the sentence representations of input, which are used to carry out the above eight probing tasks. The results are as shown in Table 10.

¹⁵ We select the SNLI model because NLI models show good generalization capacity for language representation [90, 91], which is supposed to be a strong test-bed for the evaluation.

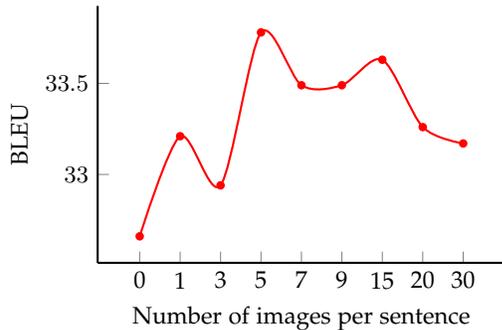


Fig. 8. Influence of the number of images on the BLEU score.

TABLE 11

Experiments on different source languages using the Multi30K dataset. The baseline is Transformer-Tiny.

Model	De-En		Fr-En	
	Test2016	Test2017	Test2016	Test2017
Baseline	42.88	40.57	54.60	48.45
UVR-TILT	43.10	40.07	54.92	49.10

Concerning semantic properties, our model gains the most significant improvement on the *SubN* and *ObjN* tasks. The result indicates that visual information helps NMT to identify and represent the subject and object information, which is consistent with our hypotheses.

6.7 Effectiveness across languages

Table 11 shows the experiment results on different source languages. We observe that our method is applicable when the source texts are in other languages such as German and French. Our proposed methods are supposed to be independent of languages because the calculation for image retrieval only relies on the light lookup table, which can be extracted from an off-the-shelf seed corpus that is available for many languages.

6.8 Joint Training and Fine-tuning

Since the multimodal and text-only machine translation tasks can benefit from the visual modality after retrieving images from the seed corpus, it is possible to bridge both tasks to train an even more powerful model. The connections between the texts and images inside the Multi30K datasets could be strong indications to bridge the gap between text and image modalities.

Therefore, we train a unified model based on UVR-CMRM by using the joint En-De datasets of Multi30K and WMT’14 (*Joint Model*), and respectively train the Multi30K (*Fine-tuned Multi30K*) and WMT’14 (*Fine-tuned WMT*) models by initializing the trainable model parameters using the joint model. According to the results shown in Table 12, we summarize the following observations:

(i) Two-stage training (joint training and fine-tuning) can boost the performance on the two concerned datasets, which indicates that the highly relevant Multi30K dataset can play the role of the seed data for training a text-only NMT model using our topic-image lookup table.

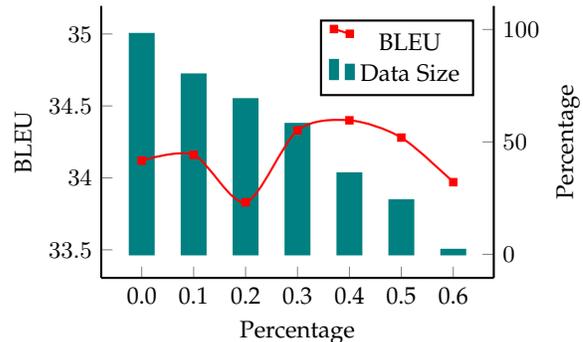


Fig. 9. BLEU score for different similarity thresholds.

TABLE 12

Results of joint training and fine-tuning.

Model	Multi30K Task	WMT Task
Joint training baseline	37.28	27.68
+ Fine-tuned Multi30K	-	27.96
+ Fine-tuned WMT	43.13	-

(ii) The major gain is achieved by fine-tuning the smaller dataset, showing that the large-scale text-only dataset with our lookup table can also provide valuable complementary information for training a multimodal model on a much smaller dataset. The result further verifies the effectiveness of our method in low-resource settings.

6.9 Parameter Sensitivity Analysis

In this section, we analyze our model sensitivity against parameter settings, including similarity threshold, number of images, and gating weight.

Influence of the similarity threshold. We investigate the influence of the similarity threshold δ that is set to filter the top-ranked images for each sentence in UVR-CMRM. Figure 9 shows the performance for thresholds in $[0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6]$ on the En-Ro test set. We observe that setting the threshold around 0.4 can yield a good balance of data size and BLEU score. It is reasonable that the best thresholds vary for different datasets because of the domain divergence of the image corpus for pre-training.

Influence of the number of images. To evaluate the influence of the number of paired images m for UVR-TILT, we constrain m in $\{0, 1, 3, 5, 7, 9, 15, 20, 30\}$ for experiments on the En-Ro test set, as shown in Figure 8. When $m = 0$, the model is the baseline NMT model, whose BLEU score is lower than all the models with images. As the number of images increases, the BLEU score also increases at the beginning (from 32.66 to 33.78) and then slightly decreases when m exceeds 5. The reason might be that too many images for a sentence would have a higher chance of noise. Therefore, we set $m = 5$ in our models.

Influence of gating weight λ . In our model, the weight λ of the gated aggregation method is learned automatically to measure the importance of the visual information. We compare by manually setting the weight λ to scalar values in $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ for experiments of UVR-TILT on the En-Ro test set. Figure 10 shows that all models with manual

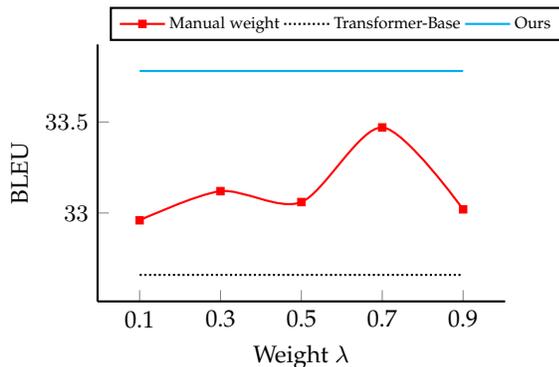


Fig. 10. Quantitative study of the gating weight λ .

λ outperform the baseline Transformer-base, indicating the effectiveness of image information. In contrast, they are inferior to the performance of our model. This means that the degree of dependency for image information varies for each source sentence, indicating the necessity of automatically learning the gating weights of image representations.

6.10 Computation Efficiency

There are mainly two extra computation costs using our method, including (i) obtaining image data for sentences and (ii) learning image representations, which are negligible compared with training an NMT model. The time of retrieving image data for MT sentences for the En-Ro dataset is less than 1 minute using GPU. The lookup table is formed as the mapping of the token (only topic words) index to the image id. Then, the retrieval method is applied as the tensor indexing from the sentence token indices (only topic words) to image ids, which is the same as the procedure of word embedding. The retrieved image ids are then sorted by frequency. Learning image representations takes about 2 minutes for all the 29,000 images in Multi30K using 6G GPU memory for feature extraction and eight CPU threads for transforming images. The extracted features are formed as the “image embedding layer” in the size of (29000, 2400) for quick access in the neural network.

7 CONCLUSIONS

This work investigates a flexible framework to incorporate visual information into sentence modeling by image retrieval from a light lookup table and learned cross-modal embedding space. Extensive empirical experiments on 14 benchmark datasets verify the effectiveness of the proposed method. A series of case studies are conducted to evaluate visual benefits and influence factors. Our method is general and can be easily implemented in existing deep-learning NLP systems for different languages. Through the proposed retrieval methods, we can provide a group of images that disclose a diversity of implicit topics that might be entailed in sentences, yielding better context grounding with fine-grained information. We show that our method enriches the representation of content words, provide fine-grained grounding information about the relationship between concepts and events, and potentially enhances the accuracy of disambiguation. Besides incorporating images to build

the pairwise relationship across modalities, it is potential to incorporate various extra knowledge as alignment topic information in the future, such as audio, not only images.

ACKNOWLEDGEMENT

Part of this study has been published as “Neural Machine Translation with Universal Visual Representation” [37] in the Eighth International Conference on Learning Representations (ICLR 2020). The extension includes three sides: (i) general tasks: this work studies the universal visual representation for language representation in a broader view of the natural language processing scenario, with experiments on 14 representative NLP tasks; (ii) new method: this work investigates new methods of semantic sentence-image matching from a shared cross-modal space, to give more accurately paired images as topic information; (iii) in-depth analysis to interpret the benefits from the visual modality.

REFERENCES

- [1] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States, 2013*, pp. 3111–3119.
- [2] J. Pennington, R. Socher, and C. Manning, “GloVe: Global vectors for word representation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543.
- [3] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018, pp. 2227–2237.
- [4] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving language understanding by generative pre-training,” *Technical report*, 2018.
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “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)*, 2019, pp. 4171–4186.
- [6] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, “Xlnet: Generalized autoregressive pretraining for language understanding,” 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, 2019*, pp. 5754–5764.
- [7] K. Zhang, G. Lv, L. Wu, E. Chen, Q. Liu, H. Wu, and F. Wu, “Image-enhanced multi-level sentence representation net for natural language inference,” in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 747–756.

- [8] M. A. McDaniel and G. O. Einstein, "Bizarre imagery as an effective memory aid: The importance of distinctiveness." *Journal of experimental psychology: Learning, memory, and cognition*, vol. 12, no. 1, p. 54, 1986.
- [9] D. Meier, *The accelerated learning handbook: A creative guide to designing and delivering faster, more effective training programs*, 2000.
- [10] J. Ive, P. Madhyastha, and L. Specia, "Distilling translations with visual awareness," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 6525–6538.
- [11] H. Shi, J. Mao, K. Gimpel, and K. Livescu, "Visually grounded neural syntax acquisition," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 1842–1861.
- [12] T. Baltruaitis, C. Ahuja, and L.-P. Morency, "Multi-modal machine learning: A survey and taxonomy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 2, pp. 423–443, 2019.
- [13] R. Hong, D. Liu, X. Mo, X. He, and H. Zhang, "Learning to compose and reason with language tree structures for visual grounding," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 684–696, 2022.
- [14] B. A. Plummer, K. J. Shih, Y. Li, K. Xu, S. Lazebnik, S. Sclaroff, and K. Saenko, "Revisiting image-language networks for open-ended phrase detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 4, pp. 2155–2167, 2022.
- [15] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, "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*.
- [16] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pre-training 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*, 2019, pp. 13–23.
- [17] H. Tan and M. Bansal, "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)*, 2019, pp. 5100–5111.
- [18] G. Li, N. Duan, Y. Fang, M. Gong, and D. Jiang, "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*, 2020, pp. 11 336–11 344.
- [19] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. Corso, and J. Gao, "Unified vision-language pre-training for image captioning and vqa," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 041–13 049.
- [20] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, "Videobert: A joint model for video and language representation learning," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, 2019, pp. 7463–7472.
- [21] E. Zablocki, B. Piwowarski, L. Soulier, and P. Gallinari, "Learning multi-modal word representation grounded in visual context," 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*, 2018, pp. 5626–5633.
- [22] Y. Meng, W. Wu, F. Wang, X. Li, P. Nie, F. Yin, M. Li, Q. Han, X. Sun, and J. Li, "Glyce: Glyph-vectors for chinese character representations," 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*, 2019, pp. 2742–2753.
- [23] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, Z. Yang, Y. Zhang, and D. Tao, "A survey on vision transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2022.
- [24] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang, "Vision-language navigation policy learning and adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 12, pp. 4205–4216, 2021.
- [25] D. Elliott, S. Frank, K. Sima'an, and L. Specia, "Multi30K: Multilingual English-German image descriptions," in *Proceedings of the 5th Workshop on Vision and Language*, 2016, pp. 70–74.
- [26] X. Li, C. Xu, X. Wang, W. Lan, Z. Jia, G. Yang, and J. Xu, "Coco-cn for cross-lingual image tagging, captioning, and retrieval," *IEEE Transactions on Multimedia*, vol. 21, no. 9, pp. 2347–2360, 2019.
- [27] Y. Yoshikawa, Y. Shigeto, and A. Takeuchi, "STAIR captions: Constructing a large-scale Japanese image caption dataset," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2017, pp. 417–421.
- [28] T. Miyazaki and N. Shimizu, "Cross-lingual image caption generation," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp. 1780–1790.
- [29] J. Hewitt, D. Ippolito, B. Callahan, R. Kriz, D. T. Wijaya, and C. Callison-Burch, "Learning translations via images with a massively multilingual image dataset," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2566–2576.
- [30] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [31] J. Gehring, M. Auli, D. Grangier, and Y. Dauphin, "A convolutional encoder model for neural machine

- translation," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2017, pp. 123–135.
- [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, 2017*, pp. 5998–6008.
- [33] M. Fadaee, A. Bisazza, and C. Monz, "Data augmentation for low-resource neural machine translation," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2017, pp. 567–573.
- [34] G. Lample, M. Ott, A. Conneau, L. Denoyer, and M. Ranzato, "Phrase-based & neural unsupervised machine translation," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 5039–5049.
- [35] X. Ma, C. Zhou, X. Li, G. Neubig, and E. Hovy, "FlowSeq: Non-autoregressive conditional sequence generation with generative flow," 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)*, 2019, pp. 4282–4292.
- [36] C. Zhou, X. Ma, J. Hu, and G. Neubig, "Handling syntactic divergence in low-resource machine translation," 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)*, 2019, pp. 1388–1394.
- [37] Z. Zhang, K. Chen, R. Wang, M. Utiyama, E. Sumita, Z. Li, and H. Zhao, "Neural machine translation with universal visual representation," in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- [38] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei *et al.*, "Oscar: Object-semantics aligned pre-training for vision-language tasks," in *European Conference on Computer Vision*. Springer, 2020, pp. 121–137.
- [39] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A deep visual-semantic embedding model," in *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States, 2013*, pp. 2121–2129.
- [40] A. Karpathy and F. Li, "Deep visual-semantic alignments for generating image descriptions," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 3128–3137.
- [41] Z. Ren, H. Jin, Z. L. Lin, C. Fang, and A. L. Yuille, "Joint image-text representation by gaussian visual-semantic embedding," in *Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016*, 2016, pp. 207–211.
- [42] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 664–676, 2017.
- [43] C. Gao, Q. Zhu, P. Wang, H. Li, Y. Liu, A. Van den Hengel, and Q. Wu, "Structured multimodal attentions for textvqa," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [44] T. Mukherjee and T. Hospedales, "Gaussian visual-linguistic embedding for zero-shot recognition," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 912–918.
- [45] L. Wang, Y. Li, J. Huang, and S. Lazebnik, "Learning two-branch neural networks for image-text matching tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 2, pp. 394–407, 2018.
- [46] M. Engilberge, L. Chevallier, P. Pérez, and M. Cord, "Finding beans in burgers: Deep semantic-visual embedding with localization," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, 2018, pp. 3984–3993.
- [47] H. Hotelling, "Relations between two sets of variates," in *Breakthroughs in statistics*, 1992, pp. 162–190.
- [48] M. Portaz, H. Randrianarivo, A. Nivaggioli, E. Maudet, C. Servan, and S. Peyronnet, "Image search using multilingual texts: a cross-modal learning approach between image and text maxime portaz qwant research," *arXiv preprint arXiv:1903.11299*, 2019.
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "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*, 2016, pp. 770–778.
- [50] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.
- [51] K. Chen, R. Wang, M. Utiyama, E. Sumita, and T. Zhao, "Neural machine translation with sentence-level topic context," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2019.
- [52] S. Xie, R. B. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 5987–5995.
- [53] T. Durand, N. Thome, and M. Cord, "WELDON: weakly supervised learning of deep convolutional neural networks," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 4743–4752.
- [54] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, 2014, pp. 1386–1393.
- [55] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 815–823.

- [56] A. Gordo, J. Almazan, J. Revaud, and D. Larlus, "End-to-end learning of deep visual representations for image retrieval," *International Journal of Computer Vision*, vol. 124, no. 2, pp. 237–254, 2017.
- [57] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- [58] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, "A large annotated corpus for learning natural language inference," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2015, pp. 632–642.
- [59] N. Nangia, A. Williams, A. Lazaridou, and S. Bowman, "The RepEval 2017 shared task: Multi-genre natural language inference with sentence representations," in *Proceedings of the 2nd Workshop on Evaluating Vector Space Representations for NLP*, 2017, pp. 1–10.
- [60] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "SQuAD: 100,000+ questions for machine comprehension of text," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 2383–2392.
- [61] L. Bentivogli, P. Clark, I. Dagan, and D. Giampiccolo, "The fifth pascal recognizing textual entailment challenge." in *ACL-PASCAL*, 2009.
- [62] W. B. Dolan and C. Brockett, "Automatically constructing a corpus of sentential paraphrases," in *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005.
- [63] Z. Chen, H. Zhang, X. Zhang, and L. Zhao, "Quora question pairs," 2018.
- [64] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia, "SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation," in *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, 2017, pp. 1–14.
- [65] A. Warstadt, A. Singh, and S. R. Bowman, "Neural network acceptability judgments," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 625–641, 2019.
- [66] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, pp. 1631–1642.
- [67] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [68] Z. Wu, L. Kong, W. Bi, X. Li, and B. Kao, "Good for misconceived reasons: An empirical revisiting on the need for visual context in multimodal machine translation," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2021, pp. 6153–6166.
- [69] Y. Yin, F. Meng, J. Su, C. Zhou, Z. Yang, J. Zhou, and J. Luo, "A novel graph-based multi-modal fusion encoder for neural machine translation," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 3025–3035.
- [70] H. Lin, F. Meng, J. Su, Y. Yin, Z. Yang, Y. Ge, J. Zhou, and J. Luo, "Dynamic context-guided capsule network for multimodal machine translation," in *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*, 2020, pp. 1320–1329.
- [71] X. Liu, P. He, W. Chen, and J. Gao, "Multi-task deep neural networks for natural language understanding," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 4487–4496.
- [72] H. Tan and M. Bansal, "Vokenization: Improving language understanding via contextualized, visually-grounded supervision," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 2066–2080.
- [73] R. Hu and A. Singh, "Unit: Multimodal multitask learning with a unified transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1439–1449.
- [74] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [75] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, "fairseq: A fast, extensible toolkit for sequence modeling," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, 2019, pp. 48–53.
- [76] P. Koehn, "Statistical significance tests for machine translation evaluation," in *Proceedings of the 2004 conference on empirical methods in natural language processing*, 2004, pp. 388–395.
- [77] J. Zhang, M. Utiyama, E. Sumita, G. Neubig, and S. Nakamura, "NICT-NAIST system for WMT17 multimodal translation task," in *Proceedings of the Second Conference on Machine Translation*, 2017, pp. 477–482.
- [78] S.-A. Grönroos, B. Huet, M. Kurimo, J. Laaksoinen, B. Meriäldo, P. Pham, M. Sjöberg, U. Sulubacak, J. Tiedemann, R. Troncy, and R. Vázquez, "The MeMAD submission to the WMT18 multimodal translation task," in *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, 2018, pp. 603–611.
- [79] O. Caglayan, P. Madhyastha, L. Specia, and L. Barrault, "Probing the need for visual context in multimodal machine translation," 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)*, 2019, pp. 4159–4170.
- [80] W. M. Brown, T. D. Gedeon, and D. I. Groves, "Use of noise to augment training data: a neural network method of mineral-potential mapping in regions of limited known deposit examples," *Natural Resources Research*, vol. 12, no. 2, pp. 141–152, 2003.
- [81] H. Noh, T. You, J. Mun, and B. Han, "Regularizing deep neural networks by noise: Its interpretation and optimization," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information*

Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, 2017, pp. 5109–5118.

- [82] J. Brownlee, “Train neural networks with noise to reduce overfitting,” *Machine Learning Mastery*, 2019.
- [83] Z. S. Harris, “Distributional structure,” *Word*, vol. 10, no. 2-3, pp. 146–162, 1954.
- [84] M. Neishi, J. Sakuma, S. Tohda, S. Ishiwatari, N. Yoshinaga, and M. Toyoda, “A bag of useful tricks for practical neural machine translation: Embedding layer initialization and large batch size,” in *Proceedings of the 4th Workshop on Asian Translation (WAT2017)*, 2017, pp. 99–109.
- [85] T. Kocmi and O. Bojar, “An exploration of word embedding initialization in deep-learning tasks,” in *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*, 2017, pp. 56–64.
- [86] S. Abnar and W. Zuidema, “Quantifying attention flow in transformers,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4190–4197.
- [87] A. Conneau, G. Kruszewski, G. Lample, L. Barrault, and M. Baroni, “What you can cram into a single $\&\#\&$ vector: Probing sentence embeddings for linguistic properties,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2126–2136.
- [88] S. Parida, O. Bojar, and S. R. Dash, “Hindi visual genome: A dataset for multimodal english-to-hindi machine translation,” *arXiv preprint arXiv:1907.08948*, 2019.
- [89] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, “Visual genome: Connecting language and vision using crowdsourced dense image annotations,” *International Journal of Computer Vision*, vol. 123, no. 1, pp. 32–73, 2017.
- [90] Z. Zhang, Y. Wu, H. Zhao, Z. Li, S. Zhang, X. Zhou, and X. Zhou, “Semantics-aware BERT for language understanding,” 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*, 2020, pp. 9628–9635.
- [91] K. Clark, M. Luong, Q. V. Le, and C. D. Manning, “ELECTRA: pre-training text encoders as discriminators rather than generators,” in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.



Zhuosheng Zhang received his Bachelor’s degree in internet of things from Wuhan University in 2016, his M.S. degree in computer science from Shanghai Jiao Tong University in 2020. He is working towards his Ph.D. degree in computer science with the Center for Brain-like Computing and Machine Intelligence of Shanghai Jiao Tong University. He was an internship research fellow at NICT from 2019-2020. His research interests include natural language processing, machine reading comprehension, dialogue systems, and

machine translation.



Kehai Chen is an Assistant Professor at Harbin Institute of Technology (Shenzhen) since 2022. Before that, he was a researcher at Japan National Institute of Information and Communications Technology (NICT) from 2018 to 2021. He received the Ph.D. degree in computer science from Harbin Institute of Technology in 2018. His research interests include machine translation and natural language processing.



Rui Wang is an associate professor at Shanghai Jiao Tong University since 2021. Before that, he was a researcher (tenured in 2020) at Japan National Institute of Information and Communications Technology (NICT) from 2016 to 2020. He received his B.S. degree from Harbin Institute of Technology in 2009, his M.S. degree from the Chinese Academy of Sciences in 2012, and his Ph.D. degree from Shanghai Jiao Tong University in 2016, all of which are in computer science. He was a joint Ph.D. at Centre National de

la Recherche Scientifique, France in 2014. His research interests are machine translation and natural language processing.

Masao Utiyama is a research manager of the National Institute of Information and Communications Technology, Japan. He completed his doctoral dissertation at the University of Tsukuba in 1997. His main research field is machine translation.



Eiichiro Sumita received the Bachelor and Master degree in computer science from The University of Electro-Communications, Japan in 1980 and 1982 and the Ph.D degree in Engineering from Kyoto University, Japan in 1999. He is currently Director of Multilingual Translation Laboratory of National Institute of Information and Communication Technology from 2006. He worked at Advanced Telecommunications Research Institute International from 1992 to 2009 and IBM Research-Tokyo from 1980 to 1991. His research interests include machine translation and e-Learning.



Zuchao Li received the B.S. degree from Wuhan University, Wuhan, China, in 2017. Since 2017, he has been a Ph.D. student with the Center for Brain-like Computing and Machine Intelligence of Shanghai Jiao Tong University, Shanghai, China. His research focuses on natural language processing, especially syntactic and semantic parsing.



Hai Zhao received the BEng degree in sensor and instrument engineering, and the MPhil degree in control theory and engineering from Yanshan University in 1999 and 2000, respectively, and the PhD degree in computer science from Shanghai Jiao Tong University, China in 2005. He is currently a full professor at department of computer science and engineering, Shanghai Jiao Tong University after he joined the university in 2009. He was a research fellow at the City University of Hong Kong from 2006 to 2009, a

visiting scholar in Microsoft Research Asia in 2011, a visiting expert in NICT, Japan in 2012. He is an ACM professional member, and served as area co-chair in ACL 2017 on Tagging, Chunking, Syntax and Parsing, (senior) area chairs in ACL 2018, 2019 on Phonology, Morphology and Word Segmentation. His research interests include natural language processing and related machine learning, data mining and artificial intelligence.