

SemVLP: Vision-Language Pre-training by Aligning Semantics at Multiple Levels

**Chenliang Li, Ming Yan, Haiyang Xu
Fuli Luo, Wei Wang, Bin Bi, Songfang Huang**

Alibaba Group

{lcl193798, ym119608, shuofeng.xhy}@alibaba-inc.com
{lf1259702, hebian.ww, b.bi, songfang.hsf}@alibaba-inc.com

Abstract

Vision-language pre-training (VLP) on large-scale image-text pairs has recently witnessed rapid progress for learning cross-modal representations. Existing pre-training methods either directly concatenate image representation and text representation at a feature level as input to a single-stream Transformer, or use a two-stream cross-modal Transformer to align the image-text representation at a high-level semantic space. In real-world image-text data, we observe that it is easy for some of the image-text pairs to align simple semantics on both modalities, while others may be related after higher-level abstraction. Therefore, in this paper, we propose a new pre-training method SemVLP, which jointly aligns both the low-level and high-level semantics between image and text representations. The model is pre-trained iteratively with two prevalent fashions: single-stream pre-training to align at a fine-grained feature level and two-stream pre-training to align high-level semantics, by employing a shared Transformer network with a pluggable cross-modal attention module. An extensive set of experiments have been conducted on four well-established vision-language understanding tasks to demonstrate the effectiveness of the proposed SemVLP in aligning cross-modal representations towards different semantic granularities.

1 Introduction

Inspired by recent development of pre-trained language models in various NLP tasks, recent studies (Lu et al., 2019; Su et al., 2019; Tan and Bansal, 2019; Chen et al., 2019b; Li et al., 2020; Yu et al., 2020) on vision-language pre-training (VLP) have pushed the limits of a variety of Vision-and-Language (V+L) tasks, which learn the semantic alignment between the different modalities by harnessing from large-scale image-text pairs.

The semantic gap between different modalities

has always been treated as one of the most significant problems in cross-modality research. In current VLP literature, there are two mainstream architectures for bridging the cross-modal semantic gap: *single-stream architecture* and *two-stream architecture*. The former such as VL-BERT (Su et al., 2019) and UNITER (Chen et al., 2019b) assumes that the underlying semantics behind the two modalities is simple and clear, and thus simply concatenates image-region features and text features as input to a single Transformer (Vaswani et al., 2017) network for early fusion in a straightforward manner. This paradigm learns the cross-modal semantic alignment from a bottom feature level by using the self-attention mechanism. Nevertheless, the design of single-stream structure treats both modality inputs equally, leaving the inherent different peculiarity of each modality not fully exploited. In contrast, the latter like LXMERT (Tan and Bansal, 2019) and ERNIE-ViL (Yu et al., 2020) first uses separate Transformer encoders to learn high-level abstraction of image and sentence representation respectively, and then combines the two modalities together with a cross-modal Transformer. This kind of design explicitly distinguishes between different modality inputs and aligns the cross-modal representations at a higher semantic level, but is usually parameter inefficient and may ignore certain more fundamental feature-level association.

In real-world image-text data, we observe that it is easy for some of the image-text pairs to align simple semantics on both modalities, while others may be related after higher-level abstraction. As shown in Figure 1, the captions of T1 are more focused on the overview of the image with coarse-level semantics, while T2 are more detailed descriptions that emphasize on the specific parts of the images. The semantic granularity spans different levels for different captions of the same images. It is essential to explicitly consider aligning semantics at multiple levels for deeply understanding the real-world



T1: A view of two streets during the night time. T2: A man sits alone on a train platform at night. T1: Upside down picture of a building surrounded by birds. T2: An image of a building with a steeple and birds flying overhead reflected in the water. T1: A group of women playing video games together. T2: Two people using an interactive gaming system while a person observes them from a couch.

Figure 1: Examples of images with two different caption text pieces from the MS COCO caption dataset, where some captions are more fine-grained than the others that are more abstract.

image-text data.

In light of this observation, we propose a new VLP pre-training architecture SemVLP, as a fusion of single-stream and two-stream architectures, which jointly aligns the image and text representation at multiple semantic levels. We observe that both the single-stream and two-stream architectures use common Transformer module, with the main difference that the latter introduces an extra CrossAttention module to allow cross-modal alignment at higher level and the different modalities are separately encoded. To complement the advantages of different architectures, we unify the two mainstream architectures by using a shared Transformer network and a pluggable cross-modal attention module, as shown in Figure 2. To conduct more fine-grained feature-level alignment, we choose a single-stream mode and directly concatenate image-region features and text features as input to the shared Transformer network for pre-training. To enhance high-level semantic alignment, we switch to a two-stream mode by separately encoding both the image and text modalities with the same shared Transformer, where a cross-modal attention module is further added to allow cross-modal fusion at a higher semantic level. The pre-training procedure is conducted iteratively so as to align the real-world image-text data at multiple semantic levels. During the iterative pre-training phase, the shared Transformer network is forced to align the semantics at multiple levels, which enables the trained model to adapt to diverse image-text pairs. In this way, we take advantages of both mainstream architectures for cross-modal fusion, where the parameters are shared to allow for different pre-training styles that regularize with each other.

We evaluate SemVLP on a variety of representative vision-language understanding tasks, including

visual question answering, natural language visual reasoning and image-text/text-image retrieval. On all these tasks, SemVLP obtains significant improvements compared to those methods that align semantics at a single fixed level, where the proposed 12-layer SemVLP model outperforms all the previous single-stream and two-stream architectures with the same model size.

The main contributions of this work can be summarized as follows: (i) We introduce SemVLP, a simple and effective VLP method to learn generic image-text representations for V+L understanding tasks. (ii) We propose a new pre-training framework that aligns cross-modal semantics at multiple levels, which can take advantages of both single-stream and two-stream architectures. To the best of our knowledge, we are among the first to think about unifying the two mainstream architectures for better aligning the cross-modal semantics. (iii) We present extensive experiments and analysis to validate the effectiveness of the proposed SemVLP model, which can obtain superior performance with a 12-layer Transformer backbone on four V+L understanding tasks.

2 Related Work

Pre-training methods have substantially advanced the NLP field in both text understanding and text generation, such as BERT (Devlin et al., 2018), ALBERT (Lan et al., 2019), GPT (Radford et al., 2018) and T5 (Raffel et al., 2019). Inspired by language pre-training, the research community starts to pay more attention to vision-language pre-training in multi-modal scenario, and many pre-training methods have been successfully applied to V+L tasks such as visual question answering (Antol et al., 2015; Hudson and Manning, 2019b) and cross-modality retrieval (Young et al., 2014; Lin

et al., 2014). In terms of the model architecture, there are mainly two broad directions to conduct vision-language pre-training. The first line uses a single-stream transformer architecture to model both image and text representations in a unified semantic space such as VLBERT (Su et al., 2019), UNITER (Chen et al., 2019b) and OSCAR (Li et al., 2020). In contrast, the other line adopts a two-stream Transformer architecture that first encodes the image and text modalities separately, and then fuses the cross-modal representations with another Transformer network. Furthermore, some other works focus on designing different pre-training tasks to learn better cross-modal representations such as ERNIE-ViL (Yu et al., 2020) and PixelBERT (Huang et al., 2020).

In this paper, we focus on the V+L understanding tasks with VLP method. Instead of choosing only one model architecture for VL pre-training, we introduce a pioneer work of fusing both the single-stream and two-stream architectures to better align the cross-modal semantics at multiple levels.

3 SemVLP Pre-training

3.1 Model Architecture

The architecture overview of SemVLP is shown in Figure 2. Inspired by the idea of sharing the encoder and decoder in Transformers for neural machine translation (Xia et al., 2019), we base the architecture of SemVLP on a shared bidirectional Transformer encoder, where a pluggable cross-modal attention module is further used to enhance high-level semantic alignment. By sharing the model parameters and adjusting the input format, SemVLP can be flexible to switch between single-stream and two-stream pre-training architectures, with the input text and image encoded in different semantic levels. In this way, we cast both the mainstream pre-training architectures into a more compact one in that there is only one copy of parameter set, which is applicable to both the low-level and high-level semantic alignment with much less parameter cost. We iteratively pre-train on the two settings towards better understanding of the real-world image-text pairs.

3.1.1 Input Embeddings

The input to SemVLP is an image and its related sentence (e.g. caption text). Each image is represented as a sequence of objects $\{o_1, \dots, o_n\}$, and each sentence is represented as a sequence of words

$\{w_1, \dots, w_m\}$. After cross-modal fusion and alignment at multiple semantic levels, SemVLP is able to generate language representations, image representations and cross-modal representations from the image-text inputs. Given the sequence of words and objects, we first introduce the methods to embed the inputs to the feature space.

Sentence Embeddings We adopt the same method as BERT (Devlin et al., 2018), which uses WordPiece tokenizer to tokenize the input sentence into sub-word tokens. The sequence of input tokens is as $\{[CLS], w_1, \dots, w_m, [SEP]\}$, where $[CLS]$ and $[SEP]$ are special tokens in BERT. The final embedding e_i for each token is generated by combining the original word embedding, segment embedding and position embedding.

Image Embeddings We use a pre-trained object detector Faster R-CNN (Ren et al., 2015) to extract the object-level image features from the image, where each object o_j is represented as a 2048-dimensional feature vector f_j . To capture the spatial information of the object, we also encode the box-level location features for each object via a 4-dimensional vector $l_j = (\frac{x_1}{W}, \frac{y_1}{H}, \frac{x_2}{W}, \frac{y_2}{H})$, where (x_1, y_1) and (x_2, y_2) denote the coordinate of the bottom-left and top-right corner while W and H are the width and height of the input image. We concatenate f_j and l_j to form a position-sensitive object feature vector, which is further transformed into o'_j using a linear projection to ensure that it has the same vector dimension as that of word embeddings. Similar to special token $[CLS]$ in sentence embeddings, we also add a special feature $[IMG]$ to denote the representation of the entire image and add it to the beginning of the input object sequence.

3.2 Shared Transformer Encoder

Given the embeddings of the words for the sentence $\{e_i\}_{i=1}^m$ and the image regions $\{o'_j\}_{j=1}^n$, the full encoder is a stacked model with L blocks, where the l 'th block consists of a self-attention module φ_S^l , a nonlinear feed forward network φ_F^l and a pluggable cross-modal attention module φ_C^l , where superscript l represents the layer index. Both the self-attention and cross-modal attention modules are based on the multi-head attention (Vaswani et al., 2017), where the feed forward network (FFN) consists of an intermediate layer and an output layer as in BERT (Devlin et al., 2018). In both the single-stream and two-stream modes, the self-attention module φ_S^l and feed forward network φ_F^l

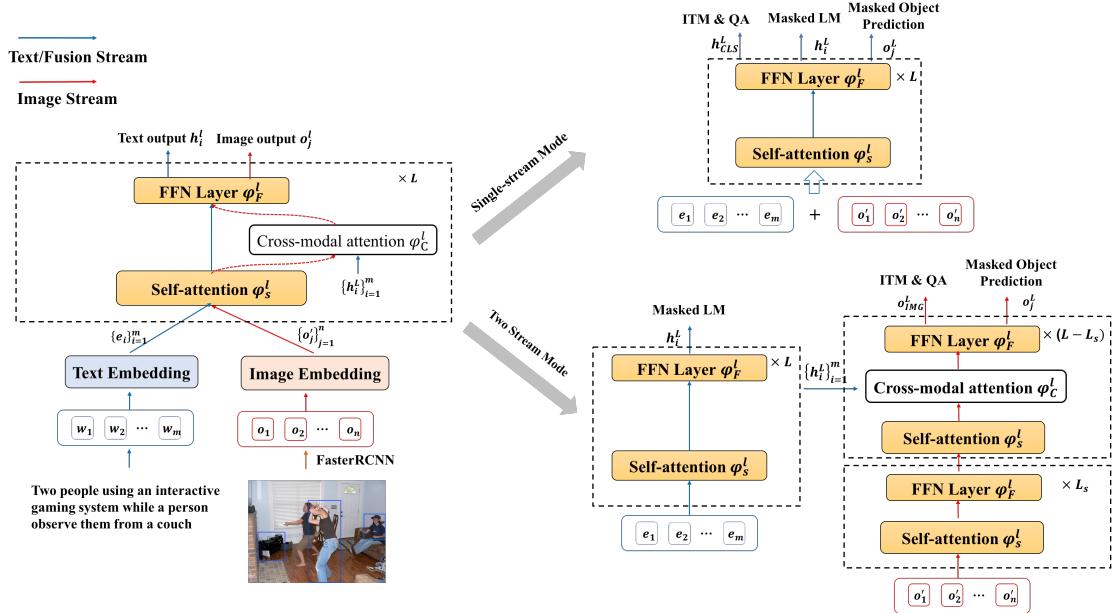


Figure 2: The overall framework of SemVLP. The whole model consists of a shared multi-layer Transformer encoder and an extra cross-modal attention module φ_C^l , where the Transformer encoder includes a self-attention module φ_S^l and a FFN layer φ_F^l (feed-forward network). In single-stream mode, we directly concatenate the image and text representation as input to the shared Transformer encoder for V+L pre-training. In two-stream mode, we reuse the same shared Transformer encoder to separately encode both the image and text representations, and add an extra cross-modal attention to align the semantics at a higher level.

are shared and tied to a single one copy of parameter set. The cross-modal attention module φ_C^l is additionally used in two-stream mode to enhance the high-level semantic alignment.

3.2.1 Feature-level Semantic Alignment

To allow a fine-grained feature-level semantic alignment, we directly concatenate the image and text embedding features as input to the single-stream mode of SemVLP, which consists of the shared self-attention module and nonlinear FFN layer. Specifically, we initialize $S^0 = \{o'_1, \dots, o'_n, e_1, \dots, e_m\}$. The encoding process can be formulated as:

$$s_i^l = \varphi_F^l(\varphi_S^l(s_i^{l-1}, S^{l-1}))$$

$$S^l = \{s_1^l, s_2^l, \dots, s_{n+m}^l\} = \{o_1^l, \dots, o_n^l, h_1^l, \dots, h_m^l\}$$

where $\{h_i^l\}$ and $\{o_j^l\}$ are the text and object representation of layer l , respectively. In this way, we can get full interaction between the image and text representations from a bottom feature-level embedding space. Eventually, we obtain $O^L = \{o_1^L, o_2^L, \dots, o_n^L\}$ and $H^L = \{h_1^L, h_2^L, \dots, h_m^L\}$, the representations of all the object outputs and text outputs of the last layer in the SemVLP encoder. The hidden representations O^L and H^L are then used to conduct the subsequent pre-training tasks.

3.2.2 High-level Semantic Alignment

For enhancing high-level semantic alignment, we adopt the two-stream mode of SemVLP, where text and image objects are separately encoded first and then fuse at a high-level semantic space. Therefore, we adopt a two-encoder architecture shown on bottom-right of Figure 2. The two-stream design of SemVLP mainly derives from the Transformer encoder-decoder network, where the main difference lies in: (1) we use both the bidirectional Transformer encoders for encoding the image and text inputs, which focuses on the V+L understanding tasks, (2) except for the cross-modal attention module, the image and text encoders share the same model parameters. We find that such parameter sharing can enhance the semantic alignment at a module level, and act as a form of regularization that stabilizes the training and saves memory consumption (Xia et al., 2019), (3) different from previous Transformer encoder-decoder architecture which introduces the cross-attention module to all blocks of the decoder, we only introduce the cross-modal attention module at the upper parts of the blocks, so as to better fuse the cross-modal representations at high-level semantic space. The encoding process of two-stream mode can be formulated

as follows:

$$\begin{aligned} h_j^l &= \varphi_F^l(\varphi_S^l(h_j^{l-1}, H^{l-1})) \\ o_j^l &= \varphi_F^l(\varphi_S^l(o_j^{l-1}, O^{l-1})), s.t. \quad l \leq L_s \\ o_{j+1}^l &= \varphi_F^l(\varphi_C^l(\varphi_S^l(o_{j+1}^{l-1}, O^{l-1}), H^L)), s.t. \quad l > L_s \end{aligned}$$

where L_s indicates the layer index that cross-modal attention is introduced. The image and text feature embeddings are first separately encoded, and then the hidden states H^L of text output in the last layer are used as input to the cross-modal attention for better understanding the image representations. Eventually, we can obtain the output representations of image objects and text, $O^L = \{o_1^L, o_2^L, \dots, o_n^L\}$ and $H^L = \{h_1^L, h_2^L, \dots, h_m^L\}$. With O^L and H^L , we could use a simple network with a softmax layer to conduct the subsequent pre-training tasks.

3.3 Joint Training

3.3.1 Pre-training Tasks

We follow LXMERT (Tan and Bansal, 2019) and use three-types of pre-training tasks: i.e., language task, vision task and cross-modality task.

Masked LM Prediction The task setup is basically the same as in BERT (Devlin et al., 2018), we randomly mask 15% tokens in the text and the model is asked to predict these masked words with the output text representations H^L . For different pre-training modes, the masked words will be predicted either with the help of visual modality so as to resolve ambiguity (single-stream mode), or from text modality alone so as to increase task difficulty (two-stream mode).

Masked Object Prediction Similarly, we pretrain the vision side by randomly masking objects, i.e., the object features are masked with zeros. We randomly mask 15% image objects and ask the model to predict properties of these masked objects with the output object representations O^L . To capture more object-level semantics, we follow the object prediction task in LXMERT (Tan and Bansal, 2019) and perform two sub-tasks: ROI-Feature Regression and Detected Label Classification. We take the detected labels output by Faster R-CNN (Ren et al., 2015) as the ground-truth labels for prediction.

Image-Text Matching (ITM) The task setup is almost the same as in LXMERT (Tan and Bansal, 2019), that we randomly sample 50% mismatched image-text pairs and 50% matched pairs, and train

an classifier to predict whether an image and a sentence match each other on the representation \mathbf{h}_{CLS}^L (single-stream mode) and \mathbf{o}_{IMG}^L (two-stream mode). One difference is that we do not enforce the masked LM prediction and Object Prediction loss when sampling a mismatched image-text pair.

Image Question Answering (QA) We also cast the image question answering task as a classification problem and pre-train the model with image QA data as in LXMERT (Tan and Bansal, 2019), which leads to a better cross-modality representation. We build the classifier on top of the representation \mathbf{h}_{CLS}^L for single-stream mode and on that of \mathbf{o}_{IMG}^L for two-stream mode.

3.3.2 Pre-training Strategy

SemVLP is pre-trained with multiple pre-training tasks and we add all these task losses with equal weights. To jointly align semantics at multiple levels, given a mini-batch of image-text pairs, 50% of the time we update the model with single-stream mode, while 50% of the time we update it with two-stream mode. In this way, for every update of SemVLP, the model is pre-trained at multiple semantic levels, so as to better model the diverse image-text data.

3.3.3 Fine-tuning Ingredient

After pre-training is completed, SemVLP can support fine-tuning with either a single-stream architecture or a two-stream architecture. In single-stream mode, the hidden state \mathbf{h}_{CLS}^L of the last layer is used for cross-modality calculation, while the hidden state of \mathbf{o}_{IMG}^L is used in two-stream mode. For each downstream task, we examine the performances for both the single-stream and two-stream fine-tuning. To yield a single model result, we use only the architecture mode with the optimal performance on development set for final evaluation.

4 Experiments

4.1 Pre-training Setup

Pre-training Data We use the same in-domain data as in LXMERT (Tan and Bansal, 2019) for pre-training. It consists of the image caption data from MS COCO (Lin et al., 2014), Visual Genome (Krishna et al., 2017), and image question answering data from VQA v2.0 (Antol et al., 2015), GQA balanced version (Hudson and Manning, 2019b) and VG-QA (Zhu et al., 2016). The total amount of the dataset is 9.18M image-and-sentence pairs on 180K

Models	Params	VQA		IR-Flickr30K			TR-Flickr30K		
		Test-dev	Test-std	R@1	R@5	R@10	R@1	R@5	R@10
Single-stream	VisualBERT	110M	70.80	71.00	-	-	-	-	-
	VLBERT	110M	71.16	-	-	-	-	-	-
	Unicoder-VL	110M	-	-	71.50	90.90	94.90	86.20	96.30
	UNITER	110M	72.70	72.91	72.52	92.36	96.08	85.90	97.10
	OSCAR	110M	73.16	73.61	-	-	-	-	-
Two-stream	ViLBERT	221M	70.55	70.92	58.20	84.90	91.52	-	-
	12-in-1	221M	73.15	-	67.90	-	-	-	-
	LXMERT	183M	72.42	72.54	-	-	-	-	-
	ERNIE-ViL	210M	72.62	72.85	74.44	92.72	95.94	86.70	97.80
Our Model	SemVLP	110M/140M	74.52	74.68	74.80	93.43	96.12	87.70	98.20

Table 1: Evaluation Results on VQA and Flickr30K.

distinct images. Besides, we also use additional out-of-domain data from Conceptual Captions (Sharma et al., 2018) and SBU Captions (Ordonez et al., 2011) for model pre-training, which consists of about 4M image-text pairs on 4M images.

Implementation Details The maximum sequence length for the sentence is set as 20. We use Faster R-CNN (Ren et al., 2015) (with ResNet-101 backbone (He et al., 2016)) pre-trained on Visual Genome dataset (Krishna et al., 2017) to detect the objects and extract the region features. We consistently keep 100 objects for each image to maximize the pre-training compute utilization by avoiding padding. For the model architecture, we pre-train a 12-layer SemVLP-base model with hidden size of 768, where we initialize it with the parameters from StructBERT base model (Wang et al., 2019). We set $L_s = 6$, which obtains the best performances on the development set of the downstream tasks, at a proper semantic level for cross-modal fusion ¹. We train SemVLP model with a total batch size of 256 for 40 epochs on 4 V100 GPUs. The Adam optimizer with initial learning rate of 1e-4 and a learning rate linear decay schedule is utilized.

4.2 Results on Downstream Tasks

We compare SemVLP model against other state-of-the-art single-stream and two-stream cross-modal pre-training models of the comparable model size ² on the following downstream tasks. The details on these tasks and fine-tuning configurations can be found in the supplementary material.

- **VQA v2.0** (Antol et al., 2015): A visual question answering task/dataset that asks a model

¹We only introduce the cross-modal attention from text space to image space due to the superior performance in our framework, where modeling of vision modality is emphasized.

²Most of the compared models have similar model size as 12-layer BERT-base.

natural language questions on a given image.

- **Image-Text Retrieval**: We test on the popular Flickr30K dataset (Young et al., 2014).
- **NLVR2** (Suhr et al., 2018): A visual reasoning task that aims to determine whether a natural language statement is true about a pair of images.
- **GQA 2019** (Hudson and Manning, 2019b): An image question answering task/dataset that emphasizes on the reasoning capability of a model to answer a question.

The results on the four downstream V+L tasks are shown in Table 1, 2, 3 respectively. We can see that: (1) Among all the VLP models of similar size to BERT base, SemVLP consistently outperforms other strong single-stream and two-stream VLP methods (e.g., UNITER (Chen et al., 2019b), OSCAR (Li et al., 2020) and 12-in-1 (Lu et al., 2020), ERNIE-ViL (Yu et al., 2020)) on all the examined tasks, which validates the effectiveness of SemVLP on combining single-stream and two-stream architectures to align semantics at multiple levels. (2) With much less parameters, the single-stream architecture can achieve comparable performance to the two-stream architecture, which is more parameter efficient. The proposed SemVLP model can be easily adapted to either architecture according to the typical scenario. By sharing parameters, SemVLP can also be parameter efficient while keeping superior performance. It is partially because SemVLP is pre-trained to align cross-modal semantics at multiple semantic levels, which makes the learning of semantic alignments more robust toward the diverse image-text pairs.

4.3 Pre-training on Different Semantic Levels

To validate the effectiveness of aligning cross-modal semantics at multiple levels, we conduct

Models	MMN (Chen et al., 2019a)	NSM (Hudson and Manning, 2019a)	LXMERT	12-in-1	OSCAR	SemVLP
Test-dev	-	-	60.00	-	61.58	62.87
Test-std	60.83	63.17	60.33	60.65	61.62	63.62

Table 2: Evaluation Results on GQA.

Models (params)	VisualBERT(110M)	LXMERT(183M)	UNITER(86M)	OSCAR(110M)	SemVLP(110M/140M)
Dev	67.40	74.90	77.14	78.07	79.00
Test-P	67.00	74.50	77.87	78.36	79.55

Table 3: Evaluation Results on NLVR2.

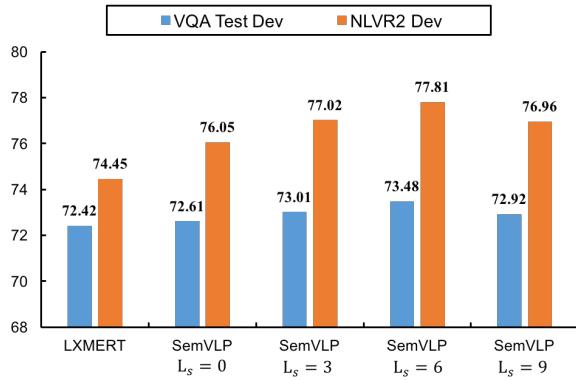


Figure 3: Results w.r.t different two-stream architectures for aligning high-level semantics on VQA and NLVR2 development set.

in-depth analysis on pre-training at different semantic levels with various architectures.

Analysis on Various Pre-training Architectures We first examine the importance of pre-training at multiple semantic levels by conducting ablation study on the pre-training fashions. Specifically, we pre-train the SemVLP model with only one type of model architecture each time and test the performance on the downstream tasks. All the pre-training settings are kept the same as in the original SemVLP pre-training. As shown in Table 4, both the two pre-training fashions play important roles in obtaining the full SemVLP model, and removing each task will consistently decrease the final downstream task performance. The single-stream architecture is used to align fair-grained feature-level semantics, while the two-stream architecture helps align semantics at a higher-level semantic space. By iterative training with a shared set of Transformer parameters, the proposed SemVLP model can take the advantage of both the single-stream architecture and two-stream architecture towards more robust vision-language pre-training.

³We only compare the single-stream and SemVLP because

Pre-training Mode	VQA	GQA	NLVR2
Baseline 1 (single-stream)	73.72	61.82	78.02
Baseline 2 (two-stream)	73.48	61.68	77.81
SemVLP	74.52	62.87	79.00

Table 4: Ablation study of different pre-training fashions on development set. Baseline 1 and Baseline 2 denote pre-training with only single-stream or only two-stream architecture, respectively.

Fine-tuning Mode	VQA	GQA	NLVR2
Single-stream	74.52	62.87	79.00
Two-stream	73.92	62.18	78.48

Table 5: Results w.r.t. different fine-tuning architectures after the SemVLP model is fully pre-trained.

Analysis on Different High-level Semantic Alignments There are many different ways for high-level semantic alignment, now we further analyze the advantage of our two-stream architecture and the specific “point” to conduct modality fusion with the cross-modal attention module. Therefore, we pre-train SemVLP with only high-level semantic alignment and examine in which layer to introduce the cross-modal attention module by setting different L_s . The pre-training details are kept the same as in the original SemVLP pre-training. We test the performance on VQA and NLVR2 tasks, and the results are shown in Figure 3. We can see that by introducing the cross-modal attention module properly, the two-stream mode of SemVLP method obtains significantly better performance than the previous two-stream model LXMERT. The best performance is obtained when $L_s = 6$, where the separated image/text encoding and cross-modal attention is equally emphasized. It again demonstrates the importance of aligning cross-modal rep-

the cross-attention module of two-stream model will interfere with the visualization.

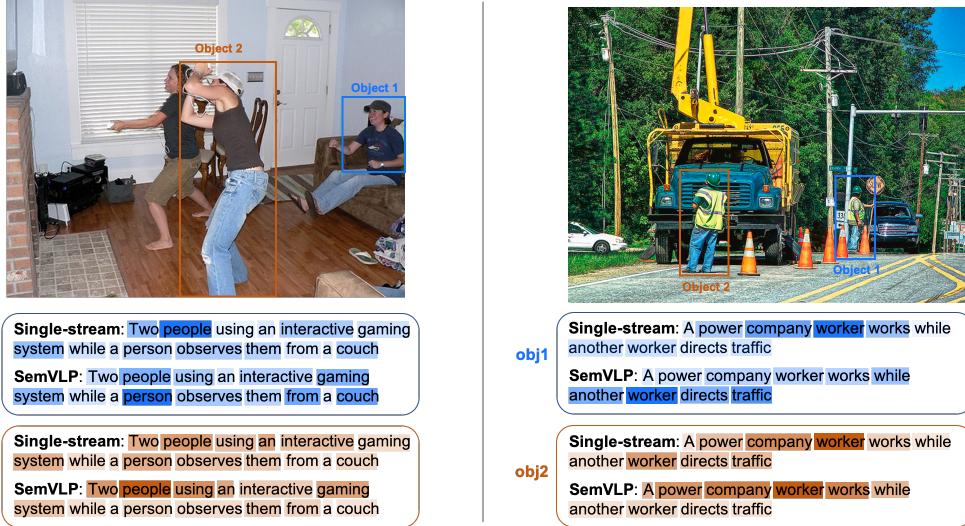


Figure 4: Visualization of the Image-to-text attention example on single-stream model and SeMVLP³. Object 1 and Object 2 are extracted by a pre-trained object detector Faster R-CNN.

resentations at a proper semantic level.

4.4 Comparison of Fine-tuning modes

After SeMVLP is fully pre-trained at multiple semantic levels, we further examine what fine-tuning architecture/mode is more appropriate for the downstream tasks. We keep the fine-tuning setting as the same with the original setting in the supplementary material, and examine the performance using different architectures. Table 5 shows that on all the examined tasks, the single-stream fine-tuning architecture gives better performance than the two-stream one does. This is due to the fact that: (1) by learning cross-modal fusion from a fine-grained feature-level, the single-stream mode can capture full association across modalities from more basic semantic granularity with powerful self-attention mechanism, which originates from the success of BERT. (2) Our design of SeMVLP shows more favor of the single Transformer encoder, which is well regularized when pre-training with both training modes via parameter sharing. The pre-training strategy tends to enhance the high-level semantic alignment into single-stream training process.

4.5 Visualization

The motivation behind SeMVLP is to align cross-modal semantics at multiple levels by taking advantages of both single-stream and two-stream architectures. As stated above, single-stream and two-stream models are good at feature-level and high-level semantic alignments, respectively. To verify this, we visualize the attention map on the same

layer and head of both the single-stream model and SeMVLP for image objects and its associated description, as shown in Figure 4. A darker color indicates a higher attention weight. Take the image on the left of Figure 4 as an example. In order to align the objects and the description, a model needs to capture the high-level text semantics: “Two people” are standing and playing games, and “a person” is sitting on the couch. For Object 1 sitting alone, the single-stream model mis-attends to “people” in the description. SeMVLP, on the other hand, properly pays high attention to “person”. For Object 2 standing alongside another, our SeMVLP attends to “people” correctly, while the single stream model fails to do so. This example shows that the single-stream cannot differentiate between the semantics of the “person” and “people” at the high level, which leads to the false alignments. The same pattern can be observed from the image on the right of Figure 4.

5 Conclusion

In this paper, we propose a new pre-training method SeMVLP to learn the joint representation of vision and language. Different from existing VLP methods relying on a fixed-level semantic alignment, we introduce to align cross-modal semantics at multiple levels, by assembling a shared Transformer encoder and a pluggable cross-modal attention module in different ways. Experiment results on various downstream V+L tasks demonstrate the effectiveness of our method for understanding the diverse semantics behind the real-world image-text data.

References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Wenhu Chen, Zhe Gan, Linjie Li, Yu Cheng, William Wang, and Jingjing Liu. 2019a. Meta module network for compositional visual reasoning. *arXiv preprint arXiv:1910.03230*.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2019b. Uniter: Universal image-text representation learning.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. 2020. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. *arXiv preprint arXiv:2004.00849*.
- Drew Hudson and Christopher D Manning. 2019a. Learning by abstraction: The neural state machine. In *Advances in Neural Information Processing Systems*, pages 5903–5916.
- Drew A Hudson and Christopher D Manning. 2019b. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6700–6709.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. *arXiv preprint arXiv:2004.06165*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23.
- Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task vision and language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10437–10446.
- Vicente Ordonez, Girish Kulkarni, and Tamara L Berg. 2011. Im2text: Describing images using 1 million captioned photographs. In *Advances in neural information processing systems*, pages 1143–1151.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. URL <https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/languageunderstandingpaper.pdf>.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*.
- Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. 2018. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.

Wei Wang, Bin Bi, Ming Yan, Chen Wu, Zuyi Bao, Jiangnan Xia, Liwei Peng, and Luo Si. 2019. Structbert: Incorporating language structures into pre-training for deep language understanding. *arXiv preprint arXiv:1908.04577*.

Yingce Xia, Tianyu He, Xu Tan, Fei Tian, Di He, and Tao Qin. 2019. Tied transformers: Neural machine translation with shared encoder and decoder. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5466–5473.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.

Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie-vil: Knowledge enhanced vision-language representations through scene graph. *arXiv preprint arXiv:2006.16934*.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004.