# Contrastive Visual-Linguistic Pretraining

**Lei Shi**[1], **Kai Shuang**[1], **Shijie Geng**[2], **Peng Su**[3], **Zhengkai Jiang**[4]
**Peng Gao**[3], **Zuohui Fu**[2], **Gerard de Melo**[5], **Sen Su**[1]
[1]Beijing University of Posts and Telecommunication [2]Rutgers University
[3]The Chinese University of Hong Kong [4]Chinese Academy of Sciences [5]Hasso Plattner Institute

## Abstract

Several multi-modality representation learning approaches such as LXMERT and ViLBERT have been proposed recently. Such approaches can achieve superior performance due to the high-level semantic information captured during large-scale multimodal pretraining. However, as ViLBERT and LXMERT adopt visual region regression and classification loss, they often suffer from domain gap and noisy label problems, based on the visual features having been pretrained on the Visual Genome dataset. To overcome these issues, we propose unbiased Contrastive Visual-Linguistic Pretraining (CVLP), which constructs a visual self-supervised loss built upon contrastive learning. We evaluate CVLP on several down-stream tasks, including VQA, GQA and NLVR2 to validate the superiority of contrastive learning on multi-modality representation learning. Our code is available at: `https://github.com/ArcherYunDong/CVLP-`.

## 1 Introduction

Language pretraining [1, 2] has revolutionized Natural Language Understanding (NLU), and strong models such as BERT [1] and RoBERTa [3] are widely used across numerous NLP tasks. Building on this, Visual-Linguistic Pretraining (VLP) has been proposed to add an extra mask-predict self-supervised strategy for the visual branch [4, 5]. Compared with VQA models that are trained from scratch such as DCN [6], BAN [7], DFAF [8] and MCAN [9], VLP relies on a similar network structure as previous methods but can achieve superior performance and better generalization ability thanks to the semantic information acquired from large-scale pretraining.

The two prominent VLP methods LXMERT [4] and ViLBERT [5] usually perform feature regression or classification for masked visual regions as the pretext task of self-supervised learning. However, we have identified several important problems: 1) Noisy Labels: $L_2$ feature regression and classification suffer from the noisy annotations in Visual Genome [10]. 2) Domain Gap: As the visual features are generated by an object detector pretrained on Visual Genome, feature regression and classification of masked regions will make the pretrained visual-linguistic model inherit the bias from Visual Genome, which results in a weak generalization ability on other downstream tasks. Taking LXMERT as an example, it can generalize better on GQA than on NLVR2 due to the overlap of pretraining and finetuning domains (the visual inputs of GQA [11] are borrowed from Visual Genome), while the images in NLVR2 [12] are collected online and consist of entirely different image manifolds compared with the sorts of images used for pretraining.

To solve the aforementioned domain gap and noisy label problems, we propose a novel Contrastive Visual-Linguistic Pretraining (CVLP), which borrows ideas from the popular contrastive learning framework in metric learning to solve the domain bias and noisy label problems. Specifically, CVLP replaces the region regression and classification with contrastive learning, which resolves the above problems. Contrastive learning aims to discriminate between positive examples and negative ones, which does not require any annotation and can solve the noisy label and domain bias problems

directly. However, due to the tremendous memory cost of Transformers [13], scaling up the batch size for contrastive learning is difficult. A conspicuous problem of contrastive learning is that the performance is highly constrained by the size of negative examples, which are bounded by the batch size. Motivated by the idea of memory banks [14, 15], we build a dynamic memory queue that caches the contextual features of the previous region and serves as negative examples in contrastive learning. The corresponding cached features drift gradually during training, thus invalidating the previously cached negative features in the memory queue. At the same time, motivated by MoCo [15], we extract features from the slowly moving query network and store them in the memory queue. When the queue is filled with features, the oldest visual contextual feature will be eliminated from the memory bank. A naive implementation of contrastive learning will fail because the network will learn to discriminate between positive and negative examples quite easily. To solve this problem, we increase feature diversity by adopting a randomly layer-dropping key network [16].

Our contributions can be summarized as below:

- We propose a novel contrastive learning framework for visual-linguistic pretraining that solves the domain bias and noisy label problems encountered with previous visual-linguistic pretraining approaches such as LXMERT and ViLBERT.

- We carry out extensive ablation studies over CVLP to validate our proposed approach. Our CVLP pretraining can achieve significant improvements over a strong baseline (LXMERT), especially when the domain gap between the pretraining and finetuning stages becomes larger. CVLP can surpass the performance of LXMERT on all three datasets (VQA, NLVR2, and GQA).



Figure 1: Example question or caption for VQA, NLVR2, GQA datasets. GQA questions are usually longer and more fine-grained than VQA ones, while NLVR2 offers a caption on a pair of images. Our CVLP consistently beats LXMERT across all three vision–language datasets.

## 2 Related Work

### 2.1 Self-supervised Learning in Vision, Language and Multi-modality

Deep Neural Networks (DNN) trained on ImageNet [17] have revolutionized automatic feature representation learning [18]. Compared to supervised training, which incurs a substantial cost for data annotation, self-supervised learning learns useful features automatically by constructing a loss from a pretext task, which does not require human annotation. In computer vision, context encoders [19] learn features by image in-painting. Jigsaw [20] learns features by predicting the position of permuted features. Kolesnikov et al. [21] carry out a large-scale study of previously proposed self-supervised learning methods and show that the performance of self-supervised tasks varies as the backbone changes. In Natural Language Understanding (NLU), large-scale pretraining with next-word prediction (GPT) [2], next sentence prediction, or masked word prediction (BERT) [1], typically trained with the Transformer architecture [13], has significantly improved the accuracy of NLU, e.g., on the GLUE benchmark [22]. Motivated by the success of self-supervised learning in both vision and language, LXMERT [4] and ViLBERT [5] have shown that masked words and visual regions can also yield a good visual-linguistic representation.

## 2.2 Contrastive Learning

Contrastive learning is a sub-branch of self-supervised learning, employing a contrastive loss to learn a representation that is useful in downstream tasks. The contrastive loss encourages the encoded instance features to be similar to positive keys while keeping away from negative ones. Different contrastive learning methods adopt different strategies to generate positive and negative keys, which is an essential factor for the quality of learned representation. [14] select the keys from a large memory bank that stores the instance features for the entire training dataset. [23, 24] generate keys using the current mini-batch samples. MoCo [15, 25] proposes a momentum encoder to generate the keys on-the-fly and store them in a fixed-size queue.

## 2.3 Multi-modality Reasoning

The backbone of current visual-linguistic pretraining is built upon previous architectures for multi-modal reasoning. Image captioning and VQA [26, 27, 28, 29] are two popular tasks that motivate the architecture design for multi-modality fusion. Specifically, attention-based architectures have widely been used in multimodal fusion. Xu et al. [30] proposed the first soft and hard attentions, showing that an attention model can yield good performance and interpretability. Yang et al. [31] proposed a multi-layer attention model by stacking attention models. Besides attention, bilinear models such as MCB [32], MLB [33] and MUTAN [34] have explored the benefit of channel interactions for multimodal understanding. Subsequently, bottom-up and top-down features [35] illustrated the benefit of employing object-level features. Recently, modeling relationships between objects and words as representation learning has been proposed in the DCN [6], BAN [7], DFAF [8], MCAN [9], QBN [36], CA-RN [37] and STSGR [38] methods.

# 3 Contrastive Visual-Linguistic Pretraining (CVLP)



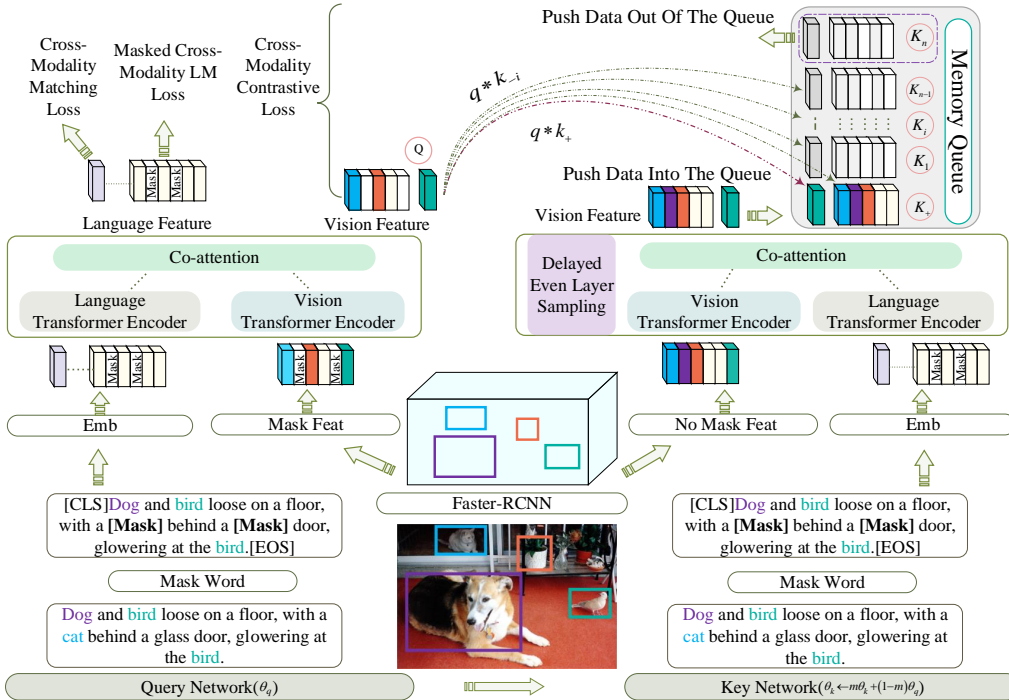Figure 2: The overall architecture of the proposed CVLP approach. CVLP includes a Query Network, a Key Network and maintains a dynamic memory queue. The entire model is trained with a combination of three cross-modality losses.

As illustrated in Figure 2, the architecture of CVLP consists of a Query Network (QueryNet) and a Key Network (KeyNet). They both contain a Language Transformer Encoder, a Vision Transformer

Encoder and a Multi-modality Fusion Transformer Encoder. At initialization, KeyNet is copied from QueryNet with the same layers and parameters. The QueryNet produces cross-modality embeddings with a masking strategy applied on both visual and textual inputs, while the KeyNet generates contextualized visual features with masking only applied to textual inputs. The output features of KeyNet are pushed into a dynamic memory queue, which continuously generates negative samples for calculating the Cross-Modality Contrastive Loss. The full CVLP model is trained with a combination of Cross-Modality Masked Language Modeling Loss, Matching Loss and Contrastive Loss. The following subsections are organized as follows: Section 3.1 introduces how visual and textual features are extracted and fused through self-attention and co-attention strategies, Sections 3.2 and 3.3 describe the design of the mask loss for the language branch and the contrastive loss for the visual branch, respectively. Section 3.4 provides further details about the dynamic memory queue mechanism and the droplayer strategy.

## 3.1 Multi-modality Fusion

Given image–sentence pairs from a vision–language dataset, we first tokenize each sentence using the WordPieces technique [39] and map a token $W_j$ to its corresponding embedding $h_{\text{emb}}(W_j) \in \mathbb{R}^{d_w}$, where $d_w = 768$. In addition, visual regions $B \in \mathbb{R}^{N \times 4}$ and their features $F \in \mathbb{R}^{N \times d_o}$ are extracted by a Faster-RNN [40] detector pretrained on Visual Genome [10] for each image $I$: $B, F = \text{RCNN}(I)$, where we detect $N = 36$ regions in each image and each region is represented using a feature dimensionality of $d_o = 2048$. Then we can calculate the visual inputs $v_i$ and textual inputs $w_j$ of CVLP as follows:

$$v_i = \frac{g_{\text{F}}(F_i) + g_{\text{P-ROI}}(B_i)}{2}, \quad w_j = h_{\text{emb}}(W_j) + h_{\text{P-word}}(P_j), \tag{1}$$

where $g_{\text{F}}$ and $g_{\text{P-ROI}}$ are two fully-connected layers that map $F_i$ and $B_i$, respectively, to the feature dimensionality $d_w$, while $h_{\text{P-word}}$ is a positional encoding function for the position $P_j$ of token $W_j$.

Taking $v_i$ and $w_j$ as inputs, CVLP adopts masking for both QueryNet and KeyNet. For QueryNet, we uniformly choose 15% of the input textual tokens for replacement. Some of the chosen tokens are replaced by the special *[MASK]* token, while the other tokens are substituted by a random token. For visual regions, we use a different masking strategy: the features of the chosen regions can either be set to zero or be replaced by region features from other images. Different from QueryNet, KeyNet only employs masking on the textual inputs, while keeping all visual region features unchanged. KeyNet and QueryNet are initialized to have the same layers and parameters. They both contain 9 layers of Language Transformer Encoder, 5 layers of Vision Transformer Encoder and 5 layers of Multi-Modality Fusion Transformer Encoder. For example, all layers in a KeyNet can be represented as:

$$\text{KeyNet} = \left\{ \begin{array}{l} r_1, r_2, r_3, r_4, r_5, \\ l_1, l_2, l_3, l_4, l_5, l_6, l_7, l_8, l_9, \\ co_1, co_2, co_3, co_4, co_5 \end{array} \right\}, \tag{2}$$

where $r_i$ stands for a self-attention layer in the visual branch, $l_i$ stands for a self-attention layer in the language branch, $co_i$ stands for a co-attention layer in the multimodality fusion branch.

The three encoders are implemented by 3 modules, namely, the visual self-attention, language self-attention and visual-linguistic co-attention modules. Visual self-attention performs information fusion between region features by using such features as both key, query and value in the attention model. We denote the key, query and value features for visual as $K_v, Q_v, V_v$, for language as $K_w, Q_w, V_w$, respectively. Then the intra-modality information fusion for visual and language features can be denoted as:

$$\widehat{v} = \text{Intra}_{\text{v} \leftarrow \text{v}}(Q_v, K_v, V_v), \quad \widehat{w} = \text{Intra}_{\text{w} \leftarrow \text{w}}(Q_w, K_w, V_w) \tag{3}$$

where the attention module of Transformer layer can be denoted as below:

$$\text{Attention}(Q, K, V) = \text{Softmax}(QK^{\text{T}}/\sqrt{d})V \tag{4}$$

After deploying intra-modality information flow for language and visual signals, we invoke an inter-modality fusion module to fuse the information from both language and visual features. The inter-modality fusion process is bi-directional, which includes information fusion from language to vision and vice versa:

$$\widetilde{v} = \text{Inter}_{\text{v} \leftarrow \text{w}}(Q_v, K_w, V_w), \quad \widetilde{w} = \text{Inter}_{\text{w} \leftarrow \text{v}}(Q_w, K_v, V_v) \tag{5}$$

4

After intra-inter modality feature fusion, we can acquire a multi-modality contextual feature embedding for each word and visual region. A contextual feature encodes the multi-modality interaction in a compact feature vector. The contextual features are used by CVLP for the mask loss in the language branch and the contrastive loss in the visual branch.

## 3.2 Mask Loss for Language Branch

In the pretraining stage, CVLP performs different pretext tasks compared with LXMERT. CVLP does not contain a supervised learning task and thus is independent of human-annotated labels. For the language branch, we keep masked language modeling and image–sentence matching prediction as two pretext tasks. Mask loss was first proposed by BERT. Subsequent visual-linguistic BERT approaches add a visual feature mask loss besides the masked language modeling loss. This loss masks out the contextual representation obtained in Section 3.1 and predicts the masked feature using its contextual information. By optimizing the mask loss, the Transformer implicitly learns to encode contextual information, which facilitates the generalization on downstream tasks. In CVLP, we only utilize mask loss for the text inputs. Additionally, we also add a matching loss, which involves a binary Yes/No classification to predict whether the sentence matches the visual feature or not. The mask loss can be formalized as follows:

$$\mathcal{L}_{\mathrm{MLM}} = -E_{w \sim D} \log P_\theta \left( w_m | \widetilde{w_{/m}} \right),\tag{6}$$

where $\theta$ denotes the parameters of the Language Transformer Encoder, $w_m$ and $\widetilde{w_{/m}}$ are the masked token to be predicted and the contextual tokens that are not masked. The matching loss is defined as:

$$\mathcal{L}_{\mathrm{MATCH}} = -E_{w_{\mathrm{CLS}} \sim D} \left[ y \log P_\theta \left( \widetilde{w_{\mathrm{CLS}}} \right) + (1 - y) \log \left[ 1 - P_\theta \left( \widetilde{w_{\mathrm{CLS}}} \right) \right] \right],\tag{7}$$

which is clearly a binary classification task. In the above equation, $w_{\mathrm{CLS}}$ stands for the *[CLS]* token which encodes the visual-linguistic information for tackling the image–sentence matching pretext task.

## 3.3 Contrastive Loss for Visual Branch

Contrastive learning performs self-supervised representation learning by discriminating visually similar feature pairs from a group of negative features. Given visual region features extracted by Faster-RCNN, we can obtain a positive query-key pair by feeding such features into both QueryNet and KeyNet. All region features from other batches are utilized as negative keys. Then we conduct contrastive learning by updating network weights to minimize the following loss:

$$\mathcal{L}_{\mathrm{CONTRAST}} = -\log \frac{\exp\left(s^+/\tau\right)}{\exp\left(s^+/\tau\right) + \sum_{j=0}^{K} \exp\left(s_j^-/\tau\right)}\tag{8}$$

$$s^+ = \widetilde{v_i^{query}} \cdot \widetilde{v_i^{key+}}, \quad s^- = \widetilde{v_i^{query}} \cdot \widetilde{v_j^{memory\_queue}}\tag{9}$$

where $\tau$ is the temperature of Softmax, $\widetilde{v_i^{key+}}$ are all positive keys of $\widetilde{v_i^{query}}$, and $\widetilde{v_j^{memory\_queue}}$ serves as negative examples for calculating $\mathcal{L}_{\mathrm{CONTRAST}}$. Traditional contrastive learning is constrained by the size of negative examples. In practice, it is time-consuming to acquire a large-sized pool of negative samples. Motivated by Momentum Contrastive (MoCo) [15], we build a dynamic visual memory queue to store the features generated by the KeyNet. The visual memory queue is empty at first, and features generated by the KeyNet are gradually placed into the queue. As training goes on, we can obtain a large visual queue to serve as negative examples. The performance of contrastive learning depends significantly on the feature diversity of the visual queue. Once the queue is full, we eliminate the oldest features. We denote the visual memory queue as:

$$\mathrm{memory\_queue} = [\widetilde{v_{b,n}^i}],\tag{10}$$

where $\widetilde{v_{b,n}^i}$ represents the visual feature that comes from the $n$-th region of the $b$-th image in the $i$-th iteration batch. One drawback of visual memory queue is feature drift during training. As the neural network is updated rapidly, the extracted features may become outdated fairly soon, which invalidates the negative examples stored in the visual queue. To resolve this, we define the weight of KeyNet as

a moving average of QueryNet when QueryNet is trained through stochastic gradient descent. The update of the network is denoted as:

$$\theta_k \leftarrow m\theta_k + (1 - m)\,\theta_q, \tag{11}$$

where $m$ stands for a momentum value, $\theta_k$ and $\theta_q$ are the parameters of KeyNet and QueryNet respectively. This form of contrastive learning can achieve superior performance due to the large visual memory queue and the small feature drift during the training progress.

### 3.4 Randomly Layer-Dropping Key Network

One important factor in training unsupervised representation learning by contrastive learning is to diversify the negative examples. Contrastive learning is highly susceptible to overfitting, thus invalidating the representation learning process. We observe that the contrastive learning loss becomes very small as the training process proceeds, suggesting that overfitting has occurred. We thus increase the diversity of features stored in the visual memory queue through the Randomly Layer-dropping Key Network. The droplayer strategy consists of a random dropout of self-attention and co-attention layers in KeyNet, which can increase the feature diversity and prevent overfitting during the training process of contrastive learning. The randomly layer-dropping Key Network can be defined as follows:

$$\text{KeyNet} = \left\{ \begin{array}{l} r_1, \text{SPL}\,(r_2)\,, r_3, \text{SPL}\,(r_4)\,, r_5, \\ l_1, \text{SPL}\,(l_2)\,, l_3, \text{SPL}\,(l_4)\,, l_5, \text{SPL}\,(l_6)\,, l_7, \text{SPL}\,(l_8)\,, l_9, \\ co_1, \text{SPL}\,(co_2)\,, co_3, \text{SPL}\,(co_4)\,, co_5, \end{array} \right\} \tag{12}$$

where SPL stands for random dropout of a layer or not. As the above equation shows, even layers in the KeyNet may be dropped during pretraining with a sampling probability of 0.5.

## 4 Experiments

In this section, we first introduce the implementation details of the proposed contrastive visual-linguistic pretraining network. Then we conduct extensive ablation studies to demonstrate the effectiveness of the proposed method. CVLP is pretrained on the same dataset as LXMERT, namely MSCOCO and Visual Genome. To assess the learned visual-linguistic features, we conduct finetuning experiments and compare CVLP with state-of-the-art methods on three downstream tasks, i.e., VQA v2 [41], GQA [11] and NLVR2 [42].
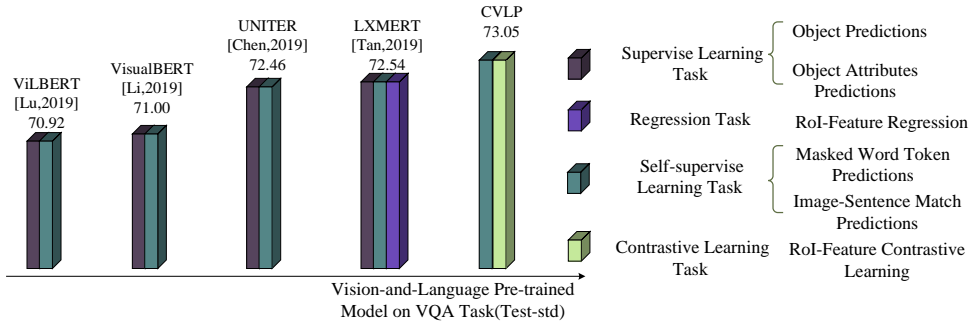


Figure 3: We show the compositions of pretext tasks used by various visual-linguistic pre-training models. Different pretext tasks require different levels of annotations and have multiple effects on downstream tasks. The height of the bars reflects the performance of each method on VQA Test-std.

**Implementation Details.** Following LXMERT, we pretrain CVLP on the same image–text pairs from MSCOCO [26] and Visual Genome [10]. In the pre-training stage, the batch size is set to 256 and the initial learning rate is 1e-4. During finetuning, the batch size is reduced to 32 and the initial learning rates of downstream tasks VQA, GQA and NLVR2 are 5e-5, 1e-5 and 5e-5, respectively. The temperature $\tau$ in the contrastive loss is set to 0.07. In both pre-training and fine-tuning stages, CVLP is optimized with Adam [43] on four Nvidia Tesla P100 GPUs.

6

## 4.1 Comparison with State-of-The-Art VLP Methods

We compare our proposed CVLP with previous visual-linguistic pretraining models, including ViLBERT [5], VisualBERT [44], UNITER [45] and LXMERT [4]. The pretraining loss utilized in each specific method is presented in Figure 3. All previous methods adopt masked visual region classification and regression. CVLP, in contrast, only needs mask loss on the text modality and contrastive learning loss on visual modality. With the help of this contrastive learning, CVLP achieves better performance on all three downstream tasks compared to previous approaches. In Table 1, we can also see that CVLP improves by 2.36% over the runner-up model UNITER on NLVR2. This improvement is the biggest among CVLP's improvements on all the three datasets, suggesting that CVLP possesses good generalization ability for large domain gap settings.

| Method | VQA | GQA | NLVR2 |
|---|---|---|---|
| | Test-dev / Test-std | Test-dev-2020 | Test-P |
| Human | - | 89.30 | 96.30 |
| Image Only | - | 17.80 | 51.90 |
| Language Only | 44.30 / - | 41.10 | 51.10 |
| LXMERT [4] | <u>72.42</u> / <u>72.54</u> | <u>61.39</u> | 74.45 |
| ViLBERT [5] | 70.55 / 70.92 | - | - |
| VisualBERT [44] (w/o Ensemble) | 70.08 / 71.00 | - | 67.00 |
| UNITER [45] | 72.27 / 72.46 | - | <u>75.58</u> |
| CVLP (finetune w/o momentum) | **72.77 / 72.90** | **61.55** | **76.20** |
| CVLP (finetune with momentum) | **72.87 / 73.05** | **61.78** | **76.81** |

Table 1: Performance comparison between CVLP and state-of-the-art visual-linguistic pretraining approaches on test splits of VQA v2, GQA 2020, and NLVR2. For all datasets, the best accuracy is in bold while the second best accuracy is underlined.

| Momentum value for pre-training | $m = 0.999$ | $m = 0.9999$ | $m = 0.99995$ | $m = 0.99999$ |
|---|---|---|---|---|
| Acc% | 70.18 | 70.40 | 70.62 | 70.08 |

Table 2: Comparison of momentum values $m$ in pre-training stage on VQA Dev-set.

| Droplayer Policy | Keynet w/o Droplayers (Epoch 1-40) | Keynet with Even Droplayers (Epoch 1-40) | Keynet with Delayed Even Droplayers (Epoch 21-40) |
|---|---|---|---|
| Acc% | 70.27 | 70.06 | 70.62 |

Table 3: Comparison of different droplayer policies on VQA Dev-set.

| Methods | VQA–Dev-set | GQA–Dev-set | NLVR2–Dev-set |
|---|---|---|---|
| No Vision Task | 66.30 | 57.10 | 50.90 |
| Feature Regression | 69.10 | 59.45 | 72.89 |
| Feature Regression + Label | 69.90 | 59.80 | 74.51 |
| Contrastive Learning | 70.62 | 59.21 | 76.47 |

Table 4: Comparison of different loss compositions on dev splits of VQA, GQA and NLVR2.

## 4.2 Ablation Studies and Analyses of CVLP

**Effects of Momentum Value.** Momentum controls the weight movement in the key network. A large momentum will result in slow drift of features. From Table 2, we can infer that a larger momentum results in a better performance on VQA because the feature drift is reduced. However, as the momentum grows to 1, the performance can drop significantly because the weight in the key

| Momentum value for fine-tuning | $m = 0.9995$ | $m = 0.99995$ | $m = 0.99997$ |
|---|---|---|---|
| NLVR2–Dev-set (1 Epoch = 2699 Iterations) | 76.47 | 72.19 | - |
| VQA–Dev-set (1 Epoch = 19753 Iterations) | 70.31 | 70.62 | - |
| GQA–Dev-set (1 Epoch = 33595 Iterations) | - | 58.98 | 59.21 |

Table 5: Comparison of momentum values $m$ in fine-tuning stage on the dev splits of NLVR2, VQA, and GQA.

network will stop to accept new information. In our experiments, we empirically determine a proper value for the momentum $m$ as $m = 1 - 1/I$, where $I$ is the iteration step in one epoch.

**Effects of Droplayer Policy.** Due to the powerful discrimination ability, contrastive learning easily overfits the training data. Our droplayer in the key network is an important technique to tackle the over-fitting problem of contrastive learning. As shown in Table 3, the key network with droplayer applied on the even layer decreases the performance. By applying a delayed droplayer policy, which takes effects after 20 epochs on even layers, the performance is significantly improved over the key network without droplayer. This experiment demonstrates the effectiveness of our proposed droplayer technique.

**Effects of Loss Composition.** In Table 4, we perform an ablation study on different loss combinations. No vision task performs visual-linguistic pretraining without adding masks on the visual features. By adding feature regression over masked visual regions, we can achieve improved performance. By adding feature regression and label classification, the results can be improved even further. After replacing feature regression and label classification loss with contrastive loss, we achieve improved performance over the three LXMERT variants on VQA and NLVR2. The performance on GQA is worse than LXMERT. This consolidates our claim that contrastive learning can perform better when the gap between pretraining and finetuning is large.
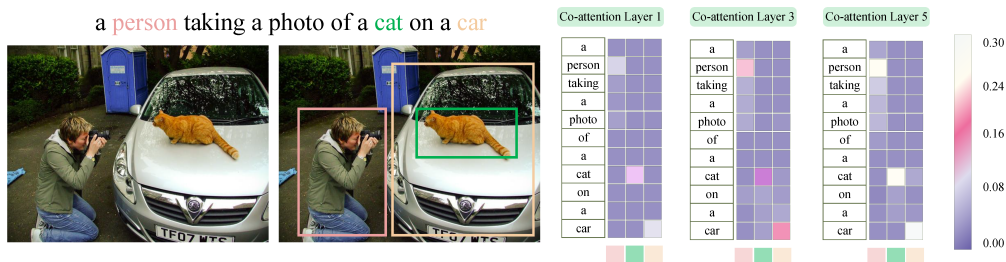


Figure 4: Illustration of attention weights in odd layers of the Co-attention Encoder. The lighter the color, the greater the weight of attention.

**Visualizing CVLP Encoder.** In Figure 4, we visualize the attention weights in odd layers (i.e., the 1st, 3rd, 5th layers) of the Co-attention Encoder in CVLP. We can see that as the layer grows, the attention weight which indicates correct word-bounding box matching also increases gradually.

# 5   Conclusion

In this paper, we propose a contrastive learning based pretraining approach for visual-linguistic representation. CVLP is not biased by the visual features pretrained on Visual Genome and can achieve superior performance, particularly when there is a substantial gap between pretraining and the downstream task. Extensive experiments and ablation studies over VQA, NLVR2, as well as GQA demonstrate the effectiveness of CVLP.

## Broader Impact

Deep learning algorithms frequently achieve superior performance on supervised tasks. However, due to their large numbers of parameters, they often require high quality and abundant training labels. Such annotation can be time-consuming and expensive. Our proposed CVLP can perform high-quality representation learning based on self-supervised pretext tasks. We believe our research can help many deep learning applications and decrease the overall cost to train and deploy a deep learning system. Large-scale pretraining with models that can cope with a domain gap have the potential to reduce possible energy usage, as one does not need to train a model from scratch for new domains. Moreover, self-supervised learning can allow us to learn from more available unlabeled data, enabling us to mitigate the well-known problems of bias and fairness in human annotations. Still, it remains important to consider the distribution of the unlabeled data to avoid biases in the model.

## References

[1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[2] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever, "Improving language understanding by generative pre-training," *URL https://s3-us-west-2. amazonaws. com/openai-assets/researchcovers/languageunsupervised/language understanding paper. pdf*, 2018.

[3] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.

[4] Hao Tan and Mohit Bansal, "Lxmert: Learning cross-modality encoder representations from transformers," *arXiv preprint arXiv:1908.07490*, 2019.

[5] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Advances in Neural Information Processing Systems*, 2019, pp. 13–23.

[6] Duy-Kien Nguyen and Takayuki Okatani, "Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6087–6096.

[7] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang, "Bilinear attention networks," in *Advances in Neural Information Processing Systems*, 2018, pp. 1564–1574.

[8] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li, "Dynamic fusion with intra-and inter-modality attention flow for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6639–6648.

[9] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian, "Deep modular co-attention networks for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6281–6290.

[10] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David 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.

[11] Drew A Hudson and Christopher D Manning, "Gqa: A new dataset for real-world visual reasoning and compositional question answering," *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[12] Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi, "A corpus of natural language for visual reasoning," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2017, pp. 217–223.

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

[14] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3733–3742.

[15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick, "Momentum contrast for unsupervised visual representation learning," *arXiv preprint arXiv:1911.05722*, 2019.

[16] Angela Fan, Edouard Grave, and Armand Joulin, "Reducing transformer depth on demand with structured dropout," *arXiv preprint arXiv:1909.11556*, 2019.

[17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[19] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2536–2544.

[20] Mehdi Noroozi and Paolo Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *European Conference on Computer Vision*. Springer, 2016, pp. 69–84.

[21] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer, "Revisiting self-supervised visual representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2019, pp. 1920–1929.

[22] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," *arXiv preprint arXiv:1804.07461*, 2018.

[23] Yonglong Tian, Dilip Krishnan, and Phillip Isola, "Contrastive multiview coding," *arXiv preprint arXiv:1906.05849*, 2019.

[24] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020.

[25] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.

[26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[27] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh, "Vqa: Visual question answering," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2425–2433.

[28] Peng Gao, Hongsheng Li, Shuang Li, Pan Lu, Yikang Li, Steven CH Hoi, and Xiaogang Wang, "Question-guided hybrid convolution for visual question answering," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 469–485.

[29] Peng Gao, Haoxuan You, Zhanpeng Zhang, Xiaogang Wang, and Hongsheng Li, "Multi-modality latent interaction network for visual question answering," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 5825–5835.

[30] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*, 2015, pp. 2048–2057.

[31] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola, "Stacked attention networks for image question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 21–29.

[32] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach, "Multimodal compact bilinear pooling for visual question answering and visual grounding," *arXiv preprint arXiv:1606.01847*, 2016.

[33] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang, "Hadamard product for low-rank bilinear pooling," *arXiv preprint arXiv:1610.04325*, 2016.

[34] Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome, "Mutan: Multimodal tucker fusion for visual question answering," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2612–2620.

[35] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6077–6086.

[36] Lei Shi, Shijie Geng, Kai Shuang, Chiori Hori, Songxiang Liu, Peng Gao, and Sen Su, "Multi-layer content interaction through quaternion product for visual question answering," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 4412–4416.

[37] Shijie Geng, Ji Zhang, Zuohui Fu, Peng Gao, Hang Zhang, and Gerard de Melo, "Character matters: Video story understanding with character-aware relations," *arXiv preprint arXiv:2005.08646*, 2020.

[38] Shijie Geng, Peng Gao, Chiori Hori, Jonathan Le Roux, and Anoop Cherian, "Spatio-temporal scene graphs for video dialog," *arXiv preprint arXiv:2007.03848*, 2020.

[39] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al., "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.

[40] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.

[41] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh, "Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[42] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi, "A corpus for reasoning about natural language grounded in photographs," *arXiv preprint arXiv:1811.00491*, 2018.

[43] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[44] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang, "Visualbert: A simple and performant baseline for vision and language," *arXiv preprint arXiv:1908.03557*, 2019.

[45] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu, "Uniter: Learning universal image-text representations," *arXiv preprint arXiv:1909.11740*, 2019.