# Retrieval-based Knowledge Augmented Vision Language Pre-training

Jiahua Rao*
School of Computer Science and
Engineering, Sun Yat-sen University
Guangzhou, China
raojh6@mail2.sysu.edu.cn

Zifei Shan*
WeChat, Tencent
Shenzhen, China
zifeishan@tencent.com

Longpo Liu
WeChat, Tencent
Shenzhen, China
longpoliu@tencent.com

Yao Zhou
WeChat, Tencent
Shenzhen, China
yaoozhou@tencent.com

Yuedong Yang[†]
School of Computer Science and
Engineering, Sun Yat-sen University
Guangzhou, China
yangyd25@mail.sysu.edu.cn

## ABSTRACT

With the recent progress in large-scale vision and language representation learning, Vision Language Pre-training (VLP) models have achieved promising improvements on various multi-modal downstream tasks. Albeit powerful, these models have not fully leveraged world knowledge to their advantage. A key challenge of knowledge-augmented VLP is the lack of clear connections between knowledge and multi-modal data. Moreover, not all knowledge present in images/texts is useful, therefore prior approaches often struggle to effectively integrate knowledge, visual, and textual information. In this study, we propose **RE**trieval-based knowledge **A**ugmented **V**ision **L**anguage (**REAVL**), a novel knowledge-augmented pre-training framework to address the above issues. For the first time, we introduce a knowledge-aware self-supervised learning scheme that efficiently establishes the correspondence between knowledge and multi-modal data, and identifies informative knowledge to improve the modeling of alignment and interactions between visual and textual modalities. By adaptively integrating informative knowledge with visual and textual information, REAVL achieves new state-of-the-art performance uniformly on knowledge-based vision-language understanding and multi-modal entity linking tasks, as well as competitive results on general vision-language tasks while only using 0.2% pre-training data of the best models. Our model shows strong sample efficiency and effective knowledge utilization.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; **Natural language processing**; **Knowledge representation and reasoning**.

*Both authors contributed equally to this research.
†Corresponding author.

## KEYWORDS

Vision-Language Pre-training, Knowledge Graph, Knowledge Retrieval, Knowledge-Augmented Model

## 1 INTRODUCTION

Recent Vision-Language Pre-training (VLP) models, such as ALBEF [21], BLIP [20] and SimVLM [35], learn multi-modal representations from large-scale image-text pairs via well-designed pre-training objectives. However, despite the strong performance of these models, they generally cannot well incorporate external world knowledge in pre-training, requiring increasingly larger networks or data to cover more facts with huge computational costs. Meanwhile, large knowledge graphs (KGs), such as Wikidata [31] and ConceptNet [28], can provide structured world knowledge to multi-modal data by representing entity descriptions and their relationships, which is implicit in vision/text but comprises complementary information [33, 38]. Therefore, leveraging structured knowledge and multi-step reasoning in KGs can really complement multi-modal learning.

However, building connections between knowledge and multi-modal data remains a stiff challenge. While previous studies such as KEPLER [33], and DRAGON [38] have incorporated knowledge into pre-training models explicitly, most of them are restricted by the knowledge behind the textual modality, neglecting a large amount of knowledge in other modalities like images. They often require entity-linking tools to build connections between knowledge and text. More importantly, these models have not fully leveraged world knowledge to their advantage because there is no ground truth indicating which knowledge is helpful for multi-modal representation. MuRAG and REVEAL [5, 13] proposed to assess the usefulness of retrieved knowledge from images for language generation modeling, but they still struggle with redundant and useless knowledge. Consequently, existing works may limit their potential to incorporate informative knowledge into multi-modal data.

To this end, we propose a **RE**trieval-based knowledge **A**ugmented **V**ision-**L**anguage Pre-training (REAVL) model, which retrieves world knowledge from KGs and leverage them to improve the multi-modal representations. REAVL has two core components: a knowledge retriever that retrieves the informative knowledge given multi-modal data, and a knowledge-augmented model that fuses multi-modal data and knowledge. Concretely, we take multi-modal data and a KG as the input data, and first retrieve the related knowledge from the KG via the knowledge retriever. Then we use a knowledge-augmented model to encode these inputs into fused representations, where the multi-modal features are fused with knowledge features through the cross-attention layer in an adaptive way.

To efficiently incorporate informative knowledge into multi-modal representations, we pretrain our model by introducing a knowledge-aware self-supervised learning scheme including four different types of modules: masked language modeling (K-MLM), masked vision modeling (K-MVM), image-text contrastive learning (K-ITC), and KG link prediction (LinkPred). Specifically, K-MLM and K-MVM are the knowledge-aware masked data modeling objectives, where the original signal is reconstructed by using its masked modal inputs and the corresponding retrieved knowledge. The masked learning signals help to reward helpful knowledge and penalize uninformative ones during knowledge retrieval. And K-ITC focuses on aligning images and text in multi-modal data. Furthermore, LinkPred makes the model use the KG structure jointly with the corresponding multi-modal data to reason about missing links in KG. We conduct ablation studies and show the effectiveness and complementarity of individual self-supervised tasks.

We benchmark REAVL on multiple vision-language benchmarks, including general, knowledge-based, and multimodal entity linking (MEL) tasks. REAVL achieves state-of-the-art performance on knowledge-based and MEL tasks. Notably, on knowledge-based VQA benchmarks, REAVL achieves new state-of-the-art results while utilizing fewer parameters and less knowledge than previous works. This demonstrates the superiority of our approach for retrieving and augmenting knowledge in multi-modal representations. On general tasks, REAVL also scores the best within models trained with a similar amount of data, competitive to models trained on billions of images while using only 0.2% of their data, showing the strong sample efficiency of our methodology.

To summarize, we make the following contributions:

- We proposed a novel retrieval-augmented pre-training framework that explicitly leverages large-scale knowledge graphs to assist vision-language pre-training.
- We are the first to introduce a knowledge-aware self-supervised learning scheme that efficiently identifies informative knowledge to improve the multi-modal modeling.
- We achieved new state-of-the-art performance uniformly on knowledge-based vision-language understanding and multi-modal entity linking tasks, as well as competitive results on general tasks while only using 0.2% pre-training data of the best models.

## 2 RELATED WORK

### 2.1 Vision-Language Pre-training Model

Vision-Language Pre-training has been developing rapidly in recent years. Most existing work focuses on modeling the interactions between images and text with transformer-based encoders[20, 21, 43]. And multiple well-designed cross-modal pre-training objectives have been proposed, for example, image-text matching, masked data modeling, and contrastive learning. They have achieved their impressive performance by scaling up the dataset with noisy image-text pairs collected from the web. For example, BLIP [20] effectively utilizes the noisy web data by bootstrapping the captions, where a captioner generates synthetic captions and a filter removes the noisy ones. Recent methods such as SimVLM [35], OFA [32], CoCa [39], and PaLi [6] unified the multi-modal pre-training model and self-supervised tasks to obtain state-of-the-art performance on various downstream tasks. These methods or frameworks benefit from their learning ability on the increasing web-collected pre-training data [3, 14], which also creates the challenge of huge computational costs, and thus complicates the optimization procedure.

Our work, by contrast, focuses on incorporating world knowledge into multi-modal representations, which have been demonstrated helpful and necessary for various downstream tasks. By introducing knowledge into multi-modal learning, our model can achieve better performance with low-resource pre-training data.

### 2.2 Knowledge Augmented Language Model

Knowledge integration is active research for improving language models. One line of research aims to add entity features into language models (LMs) [4, 33, 40]. For example, ERNIE [40] enhanced the text representations with their corresponding named entity mentions at the text token level while KB-VLP [4] augmented the knowledge from image-text pairs based on entity-linking tools and object detection models. A recent method DRAGON [38] bidirectionally fuses text and KG via a deep cross-modal model and learns joint reasoning over text and KG. These methods usually require entity linking tools to link text and KG entities before training, and the retrieved knowledge is not updated during training.

Another line of work is retrieval-augmented LMs [8, 11, 18], which retrieve relevant text from a corpus and integrate it into LMs as additional knowledge. For example, REALM [11] integrates Wikipedia passages as a knowledge base to benefit downstream knowledge-intensive tasks. However, current retrieval-augmented methods are often restricted to using text-only knowledge, neglecting an amount of knowledge in other modalities like images, which contain much information not covered by text. Compared to them, MuRAG [5] is the first retrieval model that is capable of retrieving knowledge presented in multiple modalities while REVEAL [13] construct a large-scale memory by encoding various sources of multimodal world knowledge, including Wikipedia passage and web images with alt-text captions.

To improve on the above works, we propose to effectively retrieve informative knowledge from multiple modalities and to incorporate knowledge with multi-modal representations by knowledge-aware self-supervision learning. The fundamental rationale of our method is that it improves the retrieval steps through well-designed self-supervised learning tasks, which improves the correspondence between knowledge and multi-modal data. Another distinction is that the existing retrieval-based models typically focus on adding entity- or triplet-level knowledge rather than the subgraphs around the retrieved entity to enable richer knowledge fusion.

## 2.3 Knowledge Graph representation learning

In recent years knowledge embeddings have been extensively studied by predicting missing links in graphs. Link prediction is a fundamental task in KGs [16, 42], and various works study methods to learn KG entity and relation embeddings for link prediction, such as TransE, DistMult, and RotatE [2, 23, 30, 37]. Several works additionally use textual data or pre-training LMs to help learn KG embeddings and link prediction. While these works focus on the KG-side representations, we extend the scope and use the KG-side objective (link prediction) jointly with multimodal-side objectives to train a knowledge-augmented vision-language model.

## 3 METHODS

In this section, we first briefly introduce the model architecture of our approach in Section 3.1 and Figure 1. Then we describe the proposed Knowledge Retrieval, followed by the GNN Encoder for the retrieved knowledge (Section 3.2-3.3). In Section 3.4, we show how to effectively incorporate multi-modal data with knowledge by the Knowledge-Augmented Model. Lastly, we describe the pre-training objectives with corresponding self-supervised tasks (Section 3.5).

## 3.1 Model Architecture

At the core of our method is the idea of incorporating multi-modal representations with KGs. To achieve this, We defined the image-text pairs as $\mathcal{D} = \{(I, T) | I \in \mathcal{I}, T \in \mathcal{T}\}$, and the KG is a multi-relational graph $\mathcal{G} = \{(\mathcal{V}, \mathcal{E}, \mathcal{R})\}$ where $\mathcal{V}$ is the set of entity nodes, $\mathcal{E}$ is the set of edges (triplets) and $\mathcal{R}$ is the set of relation types in the KG.

We first employ a vision encoder $f_v$ and a text encoder $f_t$ to produce the representations of the input image $I$ and text $T$. For the vision encoder, we use a 12-layer visual transformer ViT-B/16 [10]. An input image $I$ is first divided into $16 \times 16$ patches and encoded into a sequence of patch embeddings $v = \{v_1, ..., v_N\}$ where $N$ is the number of image patches. The text encoder, a 12-layer transformer, is initialized using the pre-trained $\text{BERT}_{base}$ [9] model and transforms an input text $T$ into a sequence of embeddings $t = \{t_{cls}, t_1, ..., t_{N_t}\}$, where $N_t$ is the number of text tokens and $t_{cls}$ is the encoding of the start token of a sentence.

The patch embeddings $v$ from image $I$ are then used for knowledge retrieval, which is to find the nearest entities behind the image. The set of retrieved entities is defined as $M \subset \mathcal{V}$. We then add their one-hop neighbor nodes to construct the entity subgraph $G \subset \mathcal{G}$ and feed them into the GNN Encoder. Therefore, we could obtain the retrieved entity embedding $e = \{e_1, ..., e_k\}$ after aggregating their neighbor information.

Finally, the image and text embedding as well as entity embedding are further processed by the knowledge-augmented model. The knowledge-augmented model is initialized by the first and last layer of the BERT, which enhances multi-modal representation by interacting with other modalities through cross-attention blocks.

## 3.2 Knowledge Retrieval

Instead of directly linking entity mentions in the input text to entity nodes in KGs [38], we followed the strategy of REALM [11] by explicitly asking the model to decide what knowledge to retrieve and incorporate.

In the retrieval step, the image patches with fine-grained descriptions are helpful for understanding image content and retrieving possibly helpful entities $M$ from the KG $\mathcal{G}$. Therefore, with the $N$ patches of the input image $I$, the vision encoder $f_v$ is applied to map each patch to a dense representation $v = \{v_0, v_1, ..., v_N\}$. For knowledge entries in KGs $m \in \mathcal{V}$ ($\mathcal{V}$ is the set of entity nodes), we apply the backbone encoder $f_e$ to encode the entity $m$ with its description. Formally, we define the relevance score between the image patch $v_i$ and the entity $m$ as:

$$sim(v_i, m_j) = f_v(v_i)^T f_e(m_j) \tag{1}$$

In total, we retrieve the top $N \times K$ knowledge entries relevant to image $I$. We keep top-k knowledge entries ranked by similarity scores as the set of retrieved entities $M = \{m_1, m_2, ..., m_k\}$.

In practice, we use CLIP model [25] to encode all of the entity descriptions as the candidate memory and index them using FAISS [15] to find the top-K nearest neighbors for each image patch.

For each resulting set of retrieved entities, we further add their one-hop neighbor nodes from KG to construct the entity subgraph $G$. The subgraph contains richer entities and relations, enabling our model to obtain more information about world knowledge.

## 3.3 GNN Encoder

After extracting the 2-hop entity subgraph around the image-text pair, we applied a GNN encoder to model the representation of informative entities by aggregating their neighbor information and strengthening the relational information, inspired by [26, 38]. To represent the graph, we first obtain initial node embeddings $\{e_1^{(0)}, ..., e_k^{(0)}\}$ for the retrieved entities using the entity encoder. The initial embedding of the relations in the subgraph is also initialized by their description embeddings.

Then, in each layer of the GNN, the current representation of the node embeddings $\{e_1^{(l-1)}, ..., e_k^{(l-1)}\}$ is fed into the layer to perform a round of information propagation between nodes in the graph and yield pre-fused node embeddings for each entity:

$$\{e_1^{(l)}, ..., e_k^{(l)}\} = \text{GNN}(\{e_1^{(l-1)}, ..., e_k^{(l-1)}\}) \tag{2}$$

where GNN corresponds to a variant of graph attention networks and $l$ corresponds to the layer of GNN. The GNN computes node representations for each node via message passing between neighbors on the graph:

$$e_i^{(l)} = f_n\left(\sum_{e_j \in \mathcal{N}_{e_i} \cup e_i} \alpha_{ij} \cdot f_{A_{ij}}\right) + e_i^{(l-1)} \tag{3}$$

where $f_n$ is a linear transformation, $\mathcal{N}_{e_i}$ represents the neighborhood of an entity node $e_i$, and $f_{A_{ij}}$ denotes the message passing function from its neighbors $e_j$ to $e_i$. The $\alpha_{ij}$ is an attention weight that scales the message function $f_{A_{ij}}$, which could be calculated as:

$$\beta_{ij} = \frac{f_q(e_i^{(l-1)})^T f_k(e_j^{(l-1)}, r_{ij})}{\sqrt{D}} \tag{4}$$

$$\alpha_{ij} = \frac{exp(\beta_{ij})}{\sum_{e_k \in \mathcal{N}_{e_i} \cup e_i} exp(\beta_{ik})} \tag{5}$$

**Figure 1: Illustration of the REAVL model. Given multi-modal data and a large knowledge graph (KG), we utilize a Knowledge Retriever to retrieve knowledge from KG from multi-modal data and then incorporate the retrieved knowledge in multi-modal representation learning with a Knowledge-Augmented Model. We also unify four types of knowledge-aware self-supervised tasks (i.e. Masked Language Modeling) to encourage multi-modal data and KG to mutually inform each other.**

where $f_q$ and $f_k$ are linear transformations and $r_{ij}$ is a relation embedding for the relation connecting $e_i$ and $e_j$.

The message passing function $f_{A_{ij}}$ could be computed in the following manner:

$$f_{A_{ij}} = f_m(e_j^{(l-1)}, r_{ij}) \qquad (6)$$

where $f_m$ is linear transformations and $r_{ij}$ is the same as above.

## 3.4  Knowledge Augmented Model

After using a visual transformer layer, a text transformer layer, and a GNN layer for vision, text, and entities respectively, we use a knowledge-augmented model to let the two modalities fuse the knowledge information through a 2-layer cross-modalities transformer [17, 41].

In principle, we concatenate the pre-fused embeddings including visual embeddings $v$, text embeddings $t$, and entity embeddings $e$ as the input elements in the BERT model. The input sequence always starts with a special classification element ([CLS]), then goes on with visual patch embeddings, then follows up with text token elements, and ends with the retrieved entities. A special separation element ([SEP]) is inserted between the visual and text

elements, and between the text and entity elements. Therefore, four types of input elements are involved, namely, visual, textual, entity, and special elements for disambiguating different input formats. Finally, the multi-modal features are fused with knowledge features through the cross-attention block at each layer.

Let $x = \{x_1, ..., x_p\}$ be the input elements, which are embedding features of the visual, textual, and entity elements. They are processed by a multi-layer bidirectional Transformer, where the embedding features of each element are transformed layer-by-layer in the fashion of aggregating features from the other elements with adaptive attention weights. The features of the $(l + 1)$-th layer, $x^{(l+1)}$, is computed by:

$$\widehat{h}_i^{(l+1)} = \sum_{m=1}^{M} W_m^{(l+1)} \left( \sum_{j=1}^{N} A_{i,j}^m \cdot V_m^{(l+1)} x_j^l \right) \qquad (7)$$

$$h_i^{(l+1)} = \text{LayerNorm}\left( x_i^l + \widehat{h}_i^{(l+1)} \right) \qquad (8)$$

$$\widehat{x}_i^{(l+1)} = W_2^{(l+1)} \cdot \text{GELU}\left( W_1^{(l+1)} h_i^{(l+1)} \right) \qquad (9)$$

$$x_i^{(l+1)} = \text{LayerNorm}\left( h_i^{(l+1)} + \widehat{x}_i^{(l+1)} \right) \qquad (10)$$

where $m$ in Eq.7 indexes over the attention heads, and $A_{i,j}^m = \exp[(Q_m^{(l+1)} x_i^l)^T (K_m^{(l+1)} x_j^l)]$ denotes the attention weights between elements $i$ and $j$ in the $m$-th head, and $W_m^{(l+1)}$, $Q_m^{(l+1)}$, $K_m^{(l+1)}$ and $W_1^{(l+1)}$, $W_2^{(l+1)}$ are learnable weights for $m$-th attention head and $(l+1)$-th layer.

## 3.5 Pre-training objective

To ensure that the multi-modal and KG mutually inform each other, we unify four knowledge-aware self-supervised tasks: masked language modeling, masked vision modeling, KG link prediction, and image-text contrastive learning.

### 3.5.1 Knowledge-aware Masked Language Modeling (K-MLM).
MLM is a common pre-training task used for language models. Following the strategy of REALM [11], we use Span Masking Modeling to focus on the span $s$ that may require world knowledge to predict the masked tokens. Let $\widehat{T}$ denote a masked text, and $p^{msk}(I, \widehat{T}, E)$ denote the model's predicted probability for a masked token. We adopt the strategy of [20], which minimizes a cross-entropy loss:

$$\mathcal{L}_{MLM} = \mathbb{E}_{(I,\widehat{T},E)}[H(y^{msk}, p^{msk}(I, \widehat{T}, E))] \tag{11}$$

where $y^{msk}$ is a one-hot vocabulary distribution where the ground-truth token has a probability of 1.

### 3.5.2 Knowledge-aware Masked Vision Modeling (K-MVM).
For image masking, unlike MAE [12], we aim to reconstruct the invisible patches features with visible image, text, and entity features to facilitate multi-modal information and knowledge fusion. Concretely, we propose to apply a reconstruction loss to facilitate the aggregation of multi-modal and knowledge features. Formally, the reconstruction loss is defined as:

$$\mathcal{L}_{MVM} = \mathbb{E}_{(\widehat{I},T,E)}[H(y^{msk}, p^{msk}(\widehat{I}, T, E))] \tag{12}$$

where $\widehat{I}$ denotes the masked patches and $p^{msk}(\widehat{I}, T, E))$ denotes the predicted probability for masked image patch features.

### 3.5.3 Knowledge-aware Image-Text Contrastive Learning (K-ITC).
Image-Text Contrastive Learning aims to learn better uni-modal representations before fusion. Inspired by ALBEF [21], for each image and text, we calculate the softmax-normalized image-to-text and text-to-image similarity as:

$$s^{(i2t)}(I, T) = \frac{exp(s\langle I, T_m \rangle / \tau)}{\sum_{k \neq i}^{K} exp(s\langle I, T_m \rangle)} \tag{13}$$

$$s^{(t2i)}(I, T) = \frac{exp(s\langle I_m, T \rangle / \tau)}{\sum_{k \neq i}^{K} exp(s\langle I_m, T \rangle)} \tag{14}$$

$$\mathcal{L}_{ITC} = \frac{1}{2} \mathbb{E}\left([H(y^{i2t}, s^{i2t}(I, T)) + H(y^{t2i}, s^{t2i}(I, T))]\right) \tag{15}$$

### 3.5.4 KG Link Prediction (LinkPred).
While the MLM and MVM task predicts for the multi-modal and knowledge side, link prediction holds out some edges and predicts them for the input KG. As our approach takes a joint image-text and KG entities pair as input, we expect that link prediction can encourage the model to learn to use the KG structure jointly with the textual context and visual regions to reason about missing links in the KG.

Concretely, we hold out a subset of edge triplets from the input KG $S = \{(h, r, t)\} \subset E$. Then we maps each entity node ($h$ or $t$) and relation ($r$) in the KG to a vector, $h, t, r$, and defines a scoring function $\phi_r(h, t)$ to model positive/negative triplets. Herein, we consider a KG triplet scoring function $\phi_r(h, t)$ as DistMult [37]: $\phi_r(h, t) = <h, r, t>. < \cdot, \cdot, \cdot >$ denotes the trilinear dot product. A higher $\phi$ indicates a higher chance of $(h, r, t)$ being a positive triplet (edge) instead of negative (no edge). For training, we optimize the objective:

$$\begin{aligned}\mathcal{L}_{LinkPred} = &\sum_{(h,r,t) \in S} -\log\sigma(\phi_r(h, t) + \gamma) \\ &+ \frac{1}{n} \sum_{(h',r,t')} \log\sigma(\phi_r(h', t' + \gamma))\end{aligned} \tag{16}$$

where $(h', r, t')$ are $n$ negative samples corresponding to the positive triplet $(h, r, t)$, $\gamma$ is the margin, and $\sigma$ is the sigmoid function.

## 4 EXPERIMENT SETUP

## 4.1 Data

For the image-text data, following UNITER [7], we construct our pre-training data using two web datasets (Conceptual Captions, SBU Captions) and two in-domain datasets (COCO and Visual Genome). The total number of unique images is 4.0M, and the number of image-text pairs is 5.1M. We use Wikidata5M [33], a general-domain knowledge graph designed to capture background world knowledge for the KG data. It has 4M nodes and 20M edges in total. For each entity in Wikidata5M, we use its description collected from Wikipedia to produce entity embeddings.

## 4.2 Implementation Details

Our model consists of a BERT-base with 123.7M parameters and a ViT-B/16 with 85.8M parameters. We pre-train the model for 10 epochs using a batch size of 32 on 8 NVIDIA A100 GPUs. For optimization, we use the AdamW [22] optimizer with a learning rate of 5e-5 and a weight decay of 0.02. We use DistMult [37] for KG triplet scoring for the link prediction objective, with a negative exampling of 128 triplets and a margin of $\gamma = 0.05$. To pretrain the model, we perform MLM with a token masking rate of 25%, MVM with a patch masking rate of 25%, and link prediction with an edge drop rate of 15%. Appendix A provides the details of the external knowledge memory and cost.

## 4.3 Downstream evaluation tasks

We finetune and evaluate our models on the vision-language understanding tasks (OK-VQA, AOK-VQA, VQA-v2, and SNLI-VE) [1, 24, 27, 36] and entity linking tasks (WikiDiverse and WikiPerson) [29, 34]. For the vision-language understanding task, we follow the original setting and splits used by prior works [19]. Appendix B provides the full details on these tasks and data.

**Table 1: Results for vision-language pre-training methods on popular VL benchmarks. We report accuracy for OK-VQA and SNLI-VE and vqa-score for AOK-VQA and VQA-v2. The best and second-best results are marked number and <u>number</u>, respectively. The gray number indicates that the model is trained with a significantly larger number of data than our models.**

| Task | Dataset | Method | Knowledge Sources | Results dev | test |
|---|---|---|---|---|---|
| Knowledge-based Task | OK-VQA | SUPERVISED | | | |
| | | KAT | Wikidata + Frozen GPT-3 | 53.1 | |
| | | REVIVE | Wikidata + Frozen GPT-3 | <u>56.6</u> | |
| | | PRE-TRAINING | | | |
| | | ALBEF | # image 12M | 54.7 | |
| | | BLIP | # image 129M | 55.4 | |
| | | REVEAL$_{Base}$ | CC12M + Wikidata + WIT + VQA-v2 | 55.2 | |
| | | REAVL (ours) | # image 4M + Wikidata | **57.7** | |
| | AOK-VQA | PRE-TRAINING | | | |
| | | ALBEF | # images 12M | 54.5 | - |
| | | BLIP | # images 129M | <u>56.2</u> | 50.1 |
| | | REVEAL$_{Base}$ | CC12M + Wikidata + WIT + VQA-v2 | - | 50.4 |
| | | REAVL (ours) | # images 4M + Wikidata | **58.4** | **52.7** |
| General Task | VQA-v2 | BASE DATA-SIZE | | | |
| | | VL-BERT | # images 3.3M | 71.16 | - |
| | | UNITER | # images 4M | 72.70 | 72.91 |
| | | OSCAR | # images 4M | 73.16 | 73.44 |
| | | UNIMO | # images 4M | <u>75.06</u> | <u>75.27</u> |
| | | ALBEF | # images 4M | 74.54 | 74.70 |
| | | REAVL (ours) | # images 4M + Wikidata | **77.62** | **77.79** |
| | | LARGE DATA-SIZE | | | |
| | | ALBEF | # images 12M | 75.84 | 76.04 |
| | | BLIP | # images 14M | 77.54 | 77.62 |
| | | BLIP | # images 129M | 78.25 | 78.32 |
| | | SimVLM$_{base}$ | # images 1.8B | 77.87 | 78.14 |
| | SNLI-VE | BASE DATA-SIZE | | | |
| | | UNITER | # images 4M | 79.39 | 79.38 |
| | | UNIMO | # images 4M | <u>81.11</u> | <u>80.63</u> |
| | | ALBEF | # images 4M | 80.14 | 80.30 |
| | | REAVL (ours) | # images 4M + Wikidata | **82.41** | **82.53** |
| | | LARGE DATA-SIZE | | | |
| | | ALBEF | # images 12M | 80.80 | 80.91 |
| | | SimVLM$_{base}$ | # images 1.8B | 84.20 | 84.15 |

## 5 EXPERIMENT RESULTS

### 5.1 Evaluation on V+L understanding tasks

To examine the quality of multi-modal representation learning, we first compare REAVL on four downstream V+L tasks with state-of-the-art methods, as shown in Table 1. It is clear that our method REAVL achieves state-of-the-art performance, outperforming all existing models including supervised models (KAT, REVIVE) and VLP models (ALBEF, BLIP) on the knowledge-based task. For example, on the OK-VQA dataset, REAVL achieves a 1.94% and 4.15% relative accuracy gain over the baseline REVIVE and BLIP respectively. These improvements are consistent in the AOK-VQA dataset,

demonstrating our hypothesis that large knowledge graphs can provide complementary information to multi-modal data.

Particularly, compared to the strong baseline REVEAL, which also exploits much knowledge to improve multi-modal representation, our model shows superior performance under the same size of parameters with fewer knowledge resources, demonstrating the superiority of our self-supervised learning scheme for the retrieval of informative knowledge. This is in line with our expectations as our self-supervised learning scheme could efficiently identify informative knowledge from large-scale knowledge memory and integrate them with visual and textual information.

For the general V+L tasks, our model REAVL achieves relative improvements of 3.34% on VQA-v2 and 2.36% on SNLI-VE, with

**Table 2: Entity Linking Results on WikiDiverse and WikiPerson. R@K represents the recall of the TopK retrieved entities. The best and second-best results are marked number and <u>number</u>, respectively.**

|  | WikiDiverse | | | WikiPerson | | |
|---|---|---|---|---|---|---|
|  | R@10 | R@50 | R@100 | R@1 | R@5 | R@10 |
| ViT+BERT | 61.76 | 71.30 | 73.87 | 60.56 | 72.43 | 78.72 |
| BLINK+CLIP_N_D | 66.96 | 77.18 | 80.53 | - | - | - |
| ResNet+CLIP_N | 77.64 | 81.21 | 85.69 | 68.24 | 79.83 | 82.65 |
| ResNet+CLIP_N_D | 80.64 | 84.33 | 87.56 | 73.70 | 83.47 | 84.45 |
| CLIP | <u>82.37</u> | <u>87.82</u> | **91.04** | <u>74.55</u> | <u>84.42</u> | <u>85.15</u> |
| REAVL | **83.20** | **88.42** | <u>89.59</u> | **77.58** | **85.27** | **88.38** |

**Table 3: Ablation Study on OK-VQA.**

| Training Task / Architecture | OK-VQA |
|---|---|
| w/o Knowledge Retriever | 52.71 |
| w/o GNN Aggregation | 54.33 |
| w/o Knowledge Augmented Model | 56.61 |
| Retrieval From Textual | 53.27 |
| Retrieval From Vision (Final) | **57.72** |
| MLM | 54.38 |
| MLM+MVM | 55.60 |
| MLM+MVM+ITC | 57.24 |
| MLM+MVM+ITC+LinkPred | **57.72** |

4M pre-training images. When comparing with the models that are trained with a significantly larger number of data, our model also shows competitive performances. For example, our model outperforms ALBEF and BLIP pre-trained on >10M data on both datasets. The performance supports our view that the results of the SOTA model under larger pre-training data can be achieved through the fusion of smaller data scales and knowledge. And the slightly lower performance compared to SimVLM is expected since larger data (1.8B) can provide more information than knowledge graphs.

Similar results have been shown on other general vision-language tasks such as Image Captioning, as shown in Appendix C.1. With a small amount of data, our model achieves improvements over BLIP of 4.39% on COCO Captions and 1.62% on NoCaps. When comparing with the models that are trained with a significantly larger number of data, our model also shows competitive performances.

## 5.2 Evaluation on Entity Linking tasks

REAVL is also capable of linking the visual mentions in the image to the corresponding named entity in knowledge graphs. To demonstrate the ability, we compared REAVL with the existing multi-modal entity linking models. We followed the evaluation settings used by [34] and [29] to report the comparison results with existing methods in Table 2. Considering the format of entities in KBs, we only consider the visual (image) to textual (descriptions and names) entity linking.

As we can see, ViT+BERT has achieved reasonably good performance for R@10 (i.e., 61.76 at WikiDiverse and 78.72 at WikiPerson),

which demonstrates the effectiveness of the pre-training model. Moreover, we can see that the CLIP, which is pre-trained with about 400M image-caption pairs, has achieved the strongest baseline performances, which demonstrates the effectiveness in combining both visual description and textual description. Our model has outperformed the CLIP model on 5 out of 6 metrics across the two datasets, proving that our pre-training objectives on multi-modal data such as MLM and MVM have promoted the training of the knowledge retriever. The slightly lower performance in R@100 is expected as we select candidate entities within the Top-50 for training due to the limitation of the length of the knowledge-augmented model.

## 5.3 Ablation Study

Here, we will ablate the properties of REAVL to investigate factors that influence the model performance, as we discuss below. Table 3 and Appendix C.3 show the ablation experiment results.

*5.3.1 **Incorporating Knowledge Graphs with Multi-modal data.*** The first key contribution of our model (w.r.t. existing VLP methods) is that we incorporate knowledge graphs for V+L pre-training. We find that this significantly improves the model's performance for knowledge retrieval and integration. Compared to the REAVL without the KG retriever, jointly training the KG retriever and V+L learning substantially improves the model's performance on the OK-VQA task. Furthermore, adding the proposed GNN Aggregation and the knowledge-augmented model both enhance the model performance. This is in line with our expectation as GNN Aggregation could extract more related entities and relations while the knowledge-augmented model is designed to fuse knowledge and multi-modal data in an efficient way.

*5.3.2 **Knowledge retrieval on multi-modal data.*** Another advantage of our model is its ability to retrieve knowledge from multiple modalities like images. In order to study the necessity of visual knowledge retrieval, we perform an ablation study to see the impact of knowledge retrieval. As can be seen, the performance drop of textual retrieval is more severe on the OK-VQA dataset. This is understandable because images contain more information that is not covered by text, especially in the downstream tasks. Thus, it is necessary to retrieve knowledge from images to better facilitate the fusion of knowledge and multi-modal data.

*5.3.3 **Effect of pre-training objectives.*** In order to study the contributions of different pre-training objectives, we investigated

**Table 4: Case Study on OK-VQA dataset.**

| Question | Image | Ground-Truth | Prediction | Supporting |
|---|---|---|---|---|
| (a) What is the person in the photo wearing? | | [wetsuit, suit, wet suit, trunk] | REAVL: wetsuit ✓ | **Retrieved Entities**: {boardsport, surfer, big wave surfing, **wetsuit**, dry suit, ...} |
| (b) What kind of flowers are on the table? | | [rose] | REAVL: rose ✓ | **Retrieved Entities**: {dinnerware, Duchess of Abrantes, cake plate, tea, **rose petals**} |
| (c) Which liquid seen here comes from a citrus fruit? | | [orange juice, orange] | REAVL: juice ✗ | **Retrieved Entities**: {plateau, placemat, **fruit juice**, monkey bread, **citrus juice**, ...} |
| (d) Can you guess the location where the airoplane is seen? | | [arizona, new zealand, us, tarmac] | REAVL: airport ✗ | **Retrieved Entities**: {Canadair Regional Jet, forward-swept wing, yellow line, non-towered **airport**, ...} |
| (e) What type of temperature is this? | | [cold, cold, cold, cold, cold, cold, cold, cold, pleasant, pleasant] | REAVL: cold ✓ BLIP: cool ✗ REVIVE: warm ✗ | **Retrieved Entities**: {Strawberry train, Nankai Electric Railway (Japanese railway company), rail transport in Walt Disney Parks, **cherry blossom**} |

the four pre-training objective combinations. We report their fine-tuned results on the OK-VQA dataset in Table 3. As can be seen, only with MLM, our model only achieves an accuracy of 54.38, which slightly lags behind the baseline models (ALBEF and BLIP). Adding MVM and ITC both substantially improves the pre-trained model's performance. And the proposed LinkPred further enhances the model by reasoning complex world knowledge. Finally, our model achieves a 1.94% and 4.15% relative accuracy gain over the baseline REVIVE and BLIP respectively. The ablation study demonstrates the effectiveness and complementarity of each self-supervised task.

## 5.4 Case Study

We conduct case studies on the behavior of REAVL's knowledge retrieval and knowledge augmentation, where we analyze the contribution of retrieved entities to question answering (Table 4). We can observe that our approach can accurately retrieve informative entities after training the knowledge retriever. For instance, REAVL can retrieve useful knowledge from images (e.g., wetsuit and rose petals) to generate the correct answer in example (a)(b). From the error cases, we can see that the model still generates reasonable answers for such scenarios. For example (c), we have retrieved the informative entity *citrus juice* from the image, but failed to understand the problem well. Example (d) is quite difficult to answer but our model still could generate a reasonable answer *airport*. When comparing the existing methods, we selected the case that has been reported by REVIVE. Although the question is more complex, our

model could accurately predict the answer, but existing models, BLIP and REVIVE, struggle to predict the correct answers, which can demonstrate the potential of our proposed method.

## 6 CONCLUSION

In this paper, we presented REAVL, a novel vision-language pre-training framework to incorporate world knowledge into multi-modal representations. It not only exploits the multi-modal data for better knowledge retrieval but also fuses knowledge and multi-modal data with a knowledge-augmented model. Our experiments on knowledge-based tasks show the superiority of our approach, as it outperforms existing VLP baselines with low-resource pre-training data. At the same time, the performance on entity linking tasks also shows its excellent ability in retrieving informative knowledge from massive knowledge graphs. One limitation of REAVL is that it also constitutes large computational costs due to knowledge retrieval, and an important future research will be to extend it to make it faster and apply it to larger datasets.

# REFERENCES

[1] 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*. 2425–2433.

[2] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems* 26 (2013).

[3] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. Conceptual 12M: Pushing Web-Scale Image-Text Pre-Training To Recognize Long-Tail Visual Concepts. In *CVPR*.

[4] Kezhen Chen, Qiuyuan Huang, Yonatan Bisk, Daniel McDuff, and Jianfeng Gao. 2021. Kb-vlp: Knowledge based vision and language pretraining. In *ICML, workshop*.

[5] Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, and William W Cohen. 2022. MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text. *arXiv preprint arXiv:2210.02928* (2022).

[6] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794* (2022).

[7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *ECCV*.

[8] Zhihong Chen, Guanbin Li, and Xiang Wan. 2022. Align, reason and learn: Enhancing medical vision-and-language pre-training with knowledge. In *Proceedings of the 30th ACM International Conference on Multimedia*. 5152–5161.

[9] 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* (2018).

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ICLR* (2021).

[11] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*. PMLR, 3929–3938.

[12] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16000–16009.

[13] Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-Wei Chang, Yizhou Sun, Cordelia Schmid, David A. Ross, and Alireza Fathi. 2023. REVEAL: Retrieval-Augmented Visual-Language Pre-Training with Multi-Source Multimodal Knowledge Memory. *CVPR* (2023).

[14] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*. PMLR, 4904–4916.

[15] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.

[16] Seyed Mehran Kazemi and David Poole. 2018. Simple embedding for link prediction in knowledge graphs. *Advances in neural information processing systems* 31 (2018).

[17] Yash Khare, Viraj Bagal, Minesh Mathew, Adithi Devi, U Deva Priyakumar, and CV Jawahar. 2021. Mmbert: Multimodal bert pretraining for improved medical vqa. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 1033–1036.

[18] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

[19] Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio Savarese, and Steven CH Hoi. 2022. LAVIS: A Library for Language-Vision Intelligence. *arXiv preprint arXiv:2209.09019* (2022).

[20] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *ICML*.

[21] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems* 34 (2021), 9694–9705.

[22] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).

[23] Sijie Mai, Shuangjia Zheng, Yuedong Yang, and Haifeng Hu. 2021. Communicative message passing for inductive relation reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 4294–4302.

[24] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*. 3195–3204.

[25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. PMLR, 8748–8763.

[26] Jiahua Rao, Shuangjia Zheng, Sijie Mai, and Yuedong Yang. 2022. Communicative Subgraph Representation Learning for Multi-Relational Inductive Drug-Gene Interaction Prediction. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, 3919–3925. https://doi.org/10.24963/ijcai.2022/544 Main Track.

[27] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-OKVQA: A Benchmark for Visual Question Answering using World Knowledge. *arXiv preprint arXiv:2206.01718* (2022).

[28] Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*.

[29] Wenxiang Sun, Yixing Fan, Jiafeng Guo, Ruqing Zhang, and Xueqi Cheng. 2022. Visual Named Entity Linking: A New Dataset and A Baseline. *arXiv preprint arXiv:2211.04872* (2022).

[30] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. RotatE: Knowledge Graph Embedding by Relational Rotation in Complex Space. In *International Conference on Learning Representations*. https://openreview.net/forum?id=HkgEQnRqYQ

[31] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.

[32] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*. PMLR, 23318–23340.

[33] Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021. KEPLER: A unified model for knowledge embedding and pre-trained language representation. *Transactions of the Association for Computational Linguistics* 9 (2021), 176–194.

[34] Xuwu Wang, Junfeng Tian, Min Gui, Zhixu Li, Rui Wang, Ming Yan, Lihan Chen, and Yanghua Xiao. 2022. WikiDiverse: A Multimodal Entity Linking Dataset with Diversified Contextual Topics and Entity Types. *arXiv preprint arXiv:2204.06347* (2022).

[35] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904* (2021).

[36] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. Visual Entailment: A Novel Task for Fine-grained Image Understanding. *arXiv preprint arXiv:1901.06706* (2019).

[37] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In *International Conference on Learning Representations*.

[38] Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D. Manning, Percy Liang, and Jure Leskovec. 2022. Deep Bidirectional Language-Knowledge Graph Pretraining. In *Neural Information Processing Systems (NeurIPS)*.

[39] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917* (2022).

[40] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of ACL 2019*.

[41] Shuangjia Zheng, Yongjian Li, Sheng Chen, Jun Xu, and Yuedong Yang. 2020. Predicting drug–protein interaction using quasi-visual question answering system. *Nature Machine Intelligence* 2, 2 (2020), 134–140.

[42] Shuangjia Zheng, Jiahua Rao, Ying Song, Jixian Zhang, Xianglu Xiao, Evandro Fei Fang, Yuedong Yang, and Zhangming Niu. 2021. PharmKG: a dedicated knowledge graph benchmark for bomedical data mining. *Briefings in bioinformatics* 22, 4 (2021), bbaa344.

[43] Shuangjia Zheng, Jiahua Rao, Jixian Zhang, Ethan Cohen, Chengtao Li, and Yuedong Yang. 2022. Cross-modal Graph Contrastive Learning with Cellular Images. *bioRxiv* (2022), 2022–06.

# A  KNOWLEDGE DETAILS

## A.1  Knowledge Retriever.

We initialized our visual encoder $f_v$ by CLIP. And for the retriever, we just build the MIPS index once for simplicity and do not update the entity embedding $f_e$ at the retriever. Note that we still fine-tune the visual embedding $f_v$, so we do not perform the asynchronous index refresh as in REALM fine-tuning. It is possible that refreshing the index would further improve performance. And we update our entity embedding $f_e$ at the GNN and knowledge-augmented model.

## A.2  Knowledge Extra Cost

The extra cost of the knowledge-augmented REAVL model is most from the knowledge retriever. And we use an efficient retrieval method thanks to the FAISS. During the pre-training, we took 8 days for pre-training on 8 NVIDIA A100 GPUs with a batch size of 32 and an epoch of 20, which is slightly slower than baseline methods (ALBEF and BLIP) using the same machine. As for inference, our model took 2s per sample for VQA answer generation, compared to 1.8s of the baseline method BLIP on the same machine.

# B  DATA

## B.1  V+L understanding tasks

requires a cooperative knowledge of vision and language. We fine-tune and evaluate four vision-language understanding benchmarks, including knowledge-base tasks (OK-VQA [24] and AOK-VQA [27]) and general tasks (VQAv2 [1] and SNLI-VE [36]). Specifically, we feed the embeddings of a given question (text), an image, and the retrieved entities into the knowledge-augmented model and use a 12-layer transformer decoder to generate the answer.

For the VQA-v2, OK-VQA, and AOK-VQA, we followed the prior work [19, 21] and formulate the task as a classification problem over the most frequent answers in the training set. OK-VQA and AOK-VQA is the knowledge-based VQA benchmark that requires external knowledge to answer its questions, that is, the knowledge that is not directly present in the image input.

## B.2  Multi-modal Entity Linking

aims at linking mentions with multi-modal contexts to the referent entities from a knowledge base (e.g., Wikipedia), is an essential task for many multi-modal applications. We consider two datasets for multi-modal Entity Linking: WikiDiverse [34] and WikiPerson [29]. WikiDiverse is a high-quality human-annotated MEL dataset with diversified contextual topics and entity types from Wikinews, while WikiPerson is a high-quality human-annotated visual person-linking dataset from VisualNews. Herein, we finetuned the dual encoder of image and entity with the contrastive loss function of mention-entity samples to demonstrate the superiority of our retrieval process.

# C  ADDITIONAL EXPERIMENTS

## C.1  Experiments on Image Captioning Task.

We also evaluated our model on other general vision-language tasks such as Image Captioning. The results are consistent with the VQA and VE datasets in trend, as shown in Table 6. With a small amount

of data, our model REAVL achieves improvements over BLIP of 4.39% on COCO Captions and 1.62% on NoCaps. When comparing with the models that are trained with a significantly larger number of data, our model also shows competitive performances.

## C.2  Zero-shot Learning on entity linking task.

We add the experiments of the zero-shot tasks on the WikiPerson dataset to demonstrate the ability of the retriever, as shown in Table 7. The results demonstrate the zero-shot ability of our model and the improvement over CLIP also illustrates that our knowledge retriever has been learning in the correct direction.

## C.3  Ablation study

*C.3.1  Ablation study on the number of retrieved knowledge entries.* We have conducted the ablation study on the impact of the number of retrieved knowledge entries (Top-K). As the number of retrieved knowledge entries increases, the performance of REAVL increases rapidly thanks to the more informative knowledge. Subsequently, the performance slowly descends as K continues to increase, indicating that an overlarge of retrieved knowledge entries will bring redundant and harmful information.

*C.3.2  Analysis on Knowledge Retriever.* As shown in Table 8, we analyzed the prediction results on OK-VQA to determine whether they were caused by the error of the retriever. For instance, we compared the CLIP model with our knowledge retriever when answering the question "What is the person in the photo wearing?". CLIP model could recognize multiple objects in the image, such as waterist, wave energy, and even a company in Nepal. But the retrieved knowledge is incorrect as it does not relate to the question text, resulting in the wrong answer. In contrast, our knowledge retriever, with the help of the proposed knowledge-aware task, makes an accurate prediction by focusing on the person in the image wearing a "wetsuit" or "dry suit". It again verified the importance of our proposed knowledge-aware self-supervised tasks.

For the GNN Aggregation and Knowledge Augmented module, as shown in example (e) in Table 4, to answer the question "What type of temperature is this?", the knowledge from the input image and question is not sufficient for answering. And our methods could retrieve the entity "cherry blossom" from an image and predict the correct answer "warm" based on the climate in which cherry blossoms usually bloom. Interestingly, the useful knowledge does not directly come from the retrieved entity but rather from its neighboring knowledge ("Blooming season"). In contrast, the baseline method REVIVE without the GNN Aggregation module cannot predict the correct answer. It again verified the importance of our proposed GNN Aggregation and Knowledge Augmented module.

*C.3.3  Model Scalability.* We have demonstrated that our model can be scaled to more data. Table 9 shows the performance comparison of different data sizes. The significant improvement demonstrates the scalability of our model. We also compared the ALBEF model with our KG modeling using the same 1.3M (in-domain, COCO+VG) dataset. Both models were trained with a batch size of 32 using the same number of machines. The 8% improvements presented in the table below demonstrate the superiority of our KG modeling in the same data regime.

**Table 5: Knowledge Extra Cost**

|  | # INSTANCE | # TRIPLETS | STORAGE (INITIAL-EMBEDDING) | STORAGE (FAISS-INDEX) |
|---|---|---|---|---|
| WIKIDATA5M | 4,594,485 | 20,614,279 | 8.76GB | 2.19GB |

**Table 6: Results for vision-language pre-training methods on popular image captioning benchmarks. We report CIDEr for COCO Captions and NoCaps. The best and second-best results are marked number and <u>number</u>, respectively. The gray number indicates that the model is trained with a significantly larger number of data than our models.**

| MODEL | KNOWLEDGE RESOURCES | COCO CAPTION | NOCAPS |
|---|---|---|---|
| **BASE DATA-SIZE** | | | |
| BLIP | # IMAGE 14M | 129.7 | 105.1 |
| REVAL | # IMAGE 4M + WIKIDATA5M | 135.4 | 106.8 |
| **LARGE DATA-SIZE** | | | |
| SIMVLM-BASE | # IMAGE 1.8B | 134.8 | 94.8 |
| SIMVLM-LARGE | # IMAGE 1.8B | 142.6 | 108.5 |
| SIMVLM-HUGE | # IMAGE 1.8B | 143.3 | 110.3 |

**Table 7: Additional Results for the entity linking task. (a) Compared with the text-based methods on WikiDiverse; (b) Zero-shot learning result on WikiPerson.**

|  | RECALL@10 | RECALL@50 | RECALL@100 |
|---|---|---|---|
| BLINK | 63.63 | 73.15 | 76.03 |
| REAVL | **83.20** | **88.42** | **89.59** |
|  | ZERO-SHOT R@1 | ZERO-SHOT R@5 | ZERO-SHOT R@10 |
| CLIP | 54.41 | 68.42 | 75.32 |
| REVAL | **55.13** | **69.24** | **76.39** |

**Table 8: Case study of Knowledge retriever.**

| Model | QUESTION | IMAGE | Retrieved Knowledge |
|---|---|---|---|
| CLIP | What is the person in the photo wearing? |  | {"wave energy", "wakeboarding resort", "Dibyashwori Hydropower Ltd. ", "Surfer Riding a Wave" "waterist", } |
| REVAL | What is the person in the photo wearing? |  | {"boardsport", "surfer", "big wave surfing", "wetsuit", "dry suit"} |

**Table 9: Model Scalability**

| MODEL | PRETRAINING-DATA-SIZE | OK-VQA |
|---|---|---|
| ALBEF | CC(IN DOMAIN 1.3M) | 40.7 |
| REAVL | CC(IN DOMAIN 1.3M) | 48.6 |
| REAVL | CC4M | **53.4** |

**Table 10: Ablation study on the number of retrieved knowledge entries (Top-K).**

| MODEL | WIKIDIVERSE R@10 |
|---|---|
| REAVL (K=10) | 78.39 |
| REAVL (K=20) | 80.57 |
| REAVL (K=50) | **83.20** |
| REAVL (K=100) | 82.96 |