

Improving Commonsense in Vision-Language Models via Knowledge Graph Riddles

Shuquan Ye¹ Yujia Xie² Dongdong Chen² Yichong Xu²
Lu Yuan² Chenguang Zhu² Jing Liao^{1*}

¹ City University of Hong Kong

² Microsoft

{shuquanye2-c, jingliao}@cityu.edu.hk

{yujiaxie, dochen, Yichong.Xu, luyuan, chezhu}@microsoft.com

Abstract

This paper focuses on analyzing and improving the commonsense ability of recent popular vision-language (VL) models. Despite the great success, we observe that existing VL-models still lack commonsense knowledge/reasoning ability (e.g., “Lemons are sour”), which is a vital component towards artificial general intelligence. Through our analysis, we find one important reason is that existing large-scale VL datasets do not contain much commonsense knowledge, which motivates us to improve the commonsense of VL-models from the data perspective. Rather than collecting a new VL training dataset, we propose a more scalable strategy, i.e., “Data Augmentation with kKnowledge graph linearization for CommonsenseE capability” (DANCE). It can be viewed as one type of data augmentation technique, which can inject commonsense knowledge into existing VL datasets on the fly during training. More specifically, we leverage the commonsense knowledge graph (e.g., ConceptNet) and create variants of text description in VL datasets via bidirectional sub-graph sequentialization. For better commonsense evaluation, we further propose the first retrieval-based commonsense diagnostic benchmark. By conducting extensive experiments on some representative VL-models, we demonstrate that our DANCE technique is able to significantly improve the commonsense ability while maintaining the performance on vanilla retrieval tasks. The code and data are available at <https://github.com/pleaseconnectwifi/DANCE>.

1. Introduction

Many vision-based problems in our daily life go beyond perception and recognition. For example, when we hear people say “It tastes sour”, we need to identify they are talking about lemons on the table instead of the chocolate cake.

*Jing Liao is the corresponding author.

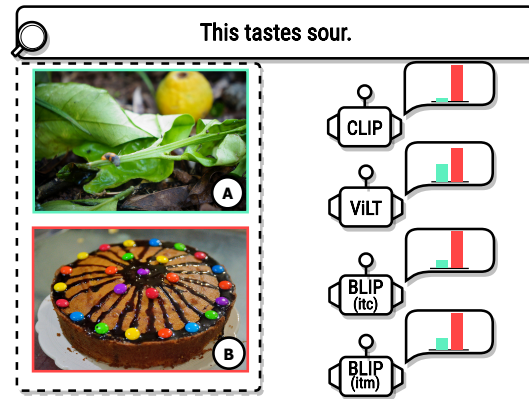


Figure 1. Illustration of the commonsense lacking problem of various popular VL-models, including CLIP [52] pre-trained with contrastive supervision, ViLT [27] with matching supervision, and BLIP [33] with the both. The bar plots suggest the alignment scores of the images to the text. All models fail in retrieving the correct image with lemon (in blue).

Therefore, it is essential for artificial general intelligence to develop commonsense capability. Vision-Language (VL) models [52] recently show promising signals on mimicking the core cognitive activities of humans by understanding the visual and textual information in the same latent space [87]. However, we observed that VL-models, e.g., CLIP [52], still struggle when minor commonsense knowledge is needed. For example, as shown in figure 1, none of the existing models correctly identify the lemon with text input “It tastes sour”.

In this work, we take a step towards injecting the VL-models with commonsense capability. More specifically, we find one important reason for the commonsense lacking issue is that existing large-scale VL datasets do not contain much commonsense knowledge. On the one hand, regular VL datasets, e.g., COCO [37] and CC 12M [9] contain more nouns and descriptive adjectives, with much fewer verbs and particles compared to regular texts. This distribution difference suggests that it might be infeasible for VL-

models to gain commonsense capability by purely enlarging the data size, unlike language-only models [24, 51]. Also, other objectives like visual question answering or generation are not widely applicable for training and have limited data size.

Inspired by the aforementioned findings, we propose Data Augmentation with kNowledge graph linearization for Commonsense capability (DANCE). The main idea is to generate commonsense-augmented image-text pairs. To do so, one natural idea is to leverage the rich commonsense knowledge in knowledge graphs [5, 66]. However, it is not trivial to inject the knowledge into image-text pairs. On the one hand, structured data like graphs usually require specific architectures [71, 88] to embed, which is troublesome. On the other hand, if we associate the external knowledge with the text in the training stage, we will need the external knowledge-augmentation process in the inference stage as well to avoid domain shift [61]. This is not desirable, since the corresponding knowledge is usually not available for the inference tasks. To address these challenges, we first re-organize the commonsense knowledge graph into entries with (entity, relation, entity) format, and pair them to the images that contain one of the entities. We then hide the name of entities in that image with demonstrative pronouns, e.g., “*this item*”. The generated descriptions are in textual form and therefore readily applicable for the training of most VL-models. More importantly, by forcing the model to memorize the relationships between entities in the training stage, such data augmentation is not needed in the inference stage. The data pair generation pipeline is automatic, scalable, and trustworthy, leveraging the existing consolidated commonsense knowledge base and the large and various collections of image-language supervision.

In addition, existing VL commonsense evaluations are restricted to visual question answering and generation which are not a good fit or well received in the majority of VL-models. Therefore, we propose a new diagnostic test set in a wider adaptable form, i.e., Image-Text and Text-Image Retrieval, to achieve a fair evaluation of the pre-trained VL-models. The set is upgraded by neighborhood hard-negative filtering to further ensure data quality.

The effectiveness of the proposed strategy is validated by not only our diagnostic test set, but also the most generally used visual question answering benchmark for commonsense [42]. Furthermore, we show that the commonsense capability of the models trained with DANCE can even generalize to unseen knowledge. We show the potential of the new train strategy and the test dataset via a deep study of its contents and baseline performance measurements across a variety of cutting-edge VL-models. Our main findings and contributions are summarized as follows:

1. We propose a novel commonsense-aware training strategy DANCE, which is compatible with the most

of VL-models. The inference stage needs no change.

2. We propose a new retrieval-based well-received commonsense benchmark to analyze a suite of VL-models and discover weaknesses that are not widely known: commonsense easy for humans (83%) is hard for current state-of-the-art VL-models (<42%).
3. We conduct extensive experiments to demonstrate the effectiveness of the proposed strategy and diagnostic test set. The datasets and all the code will be made publicly available.

2. Related Work

Vision-Language Contrastive Learning and Matching.

Vision-Language Contrastive Learning (VLC) and Matching (VLM), both of which aim to align vision and language, has been the fundamental tasks for Vision-Language model pre-training. They are the most commonly used objectives [10], and are used in the well-known foundation models [15, 73]: to name a few, CLIP [52], OwlVit [44], ALIGN [23], MDETR [25], Florence [85] with VLC supervision; ViLT [27], FLAVA [62], ViLBERT [41], UNITER [11], Unicoder [32], VisualBERT [35] utilize the VLM target; BLIP [33], ALBEF [34] uses them both. Many popular and large-scale image-text paired datasets [8, 9, 28, 37, 43, 48, 55, 58–60, 83, 84] are proposed on this task suitable for most common scenarios. Some of them target at specific cases like instructional video [43], 3D scene [8, 83], and fashion [55]. However, most of them are not targeted at commonsense knowledge or reasoning.

Image Text Retrieval. Image Text Retrieval (ITR) is a typical cross-modal downstream task needing retrieving an image that matches a description most and vice versa. ITR can be naturally performed by the VL-models pre-trained on VLC/VLM targets, and widely received in the majority of the VL-models [7] even in a zero-shot manner. Though, how to improve the commonsense knowledge of ITR still requires further study. CVSE [72] injects commonsense knowledge into VL-models for ITR using statistical correlations in caption corpus. However, such knowledge is constrained by the correlations of corpus and is not a perfect fit for ITR [7]. Various datasets [37, 40, 47, 77, 84] have been proposed to evaluate the ITR system. However, most of them do not question VL-models’ commonsense ability.

Commonsense in Vision-Language Datasets. Several studies on commonsense over visual-language input can be divided into two branches according to the task type. The first branch includes Visual Question Answering (VQA) datasets [22, 42, 75, 76, 83]. The model is required to give a natural language answer to a specific question-image pair, for which commonsense knowledge is required [38]. However, performing VQA evaluation automatically is not trivial for most of the VL-models, especially for dual-encoder

architecture, e.g., CLIP, as studied in [17]. Besides, as their data collection requires plenty of human effort, their image amount and language corpus of are quite small. Another branch focuses on generation or multi-choice tasks. Visual Commonsense Reasoning [86] evaluates the commonsense ability via two multi-choice tasks: question answering and answer justification. VisualCOMET [50] evaluates the commonsense ability via inferencing events and intents. However, they are collected from movie clips, and the commonsense knowledge required is focused on humans and events, which limits the variety of image and knowledge. Also, transferring from commonly used pre-training tasks to generation or multi-choice tasks itself is a challenging task [36, 45, 82]. To evaluate and improve the commonsense ability of VL-models, we take the first step towards automatic and direct applicable commonsense evaluation via ITR, along with a scalable learning strategy suitable for VLC/VLM with variable augmented commonsense.

Commonsense in NLP. Commonsense and knowledge representation has a long-period development in NLP [78, 79, 90], with plenty of famous knowledge bases [2, 5, 31, 63, 64, 66, 67] emerged. ConceptNet [66] is a popular and consolidated commonsense knowledge graph with 8 million nodes and 21 million edges collected from various sources including expert-created, crowd-sourcing, and games. While there are early explorations on general VL-model training methods with external knowledge [61], we are the first one that aims to train a general-purpose commonsense-aware VL-model, with no restrictions in the inference stage.

3. Data Augmentation Strategy with Knowledge Graph Linearization

We present our DANCE training strategy that enhances the commonsense ability of VL-models via learning with novel and scalable commonsense augmented data.

In short, we augment existing image-text pairs with knowledge graphs. Denote the set of image-text pairs as $\mathcal{D} = \{(i_k, t_k)\}_{k=1}^K$, where i_k and t_k are paired images and texts. Denote the knowledge graph as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the node set and \mathcal{E} is the edge set. Each edge $e \in \mathcal{E}$ is a tuple (v_i, r, w, v_j) , with $v_i, v_j \in \mathcal{V}$, r is the commonsense relationship pointing directionally from v_i to v_j , and w is the weight highlighting the importance of this edge. Here, v_i is denoted as the *head*, while v_j is the *tail*. For example, a directed edge from ConceptNet [66] takes the form as

(“Net”, “is used for”, 0.3, “catching fish”).

We explore an automatic way to generate paired commonsense-aware training data by augmenting commonsense knowledge into VL datasets. The automatic data construction process is shown in Fig. 2.

3.1. Extraction of Image-Entity Pair

To pair an image i_k with its corresponding knowledge, we first need to find out what entities are in it. Given the presence of the corresponding descriptive text t_k , a reliable way is to extract the linguistic entities from t_k via natural language processing toolkits¹. For example, for the image in the upper left of Fig. 2, the corresponding caption is “A cat with a box in an office”. Using standard toolkits, we can obtain the entities as “cat”, “box”, and “office”. In detail, we extract the linguistic entities $\mathcal{N}_k = \{\epsilon_{k,i}\}_{i=1}^I$, and further perform format cleaning that removes determiners or adjectives, and filtering out entities that are too general to be meaningful to get a subset $\hat{\mathcal{N}}_k \subseteq \mathcal{N}_k$. In this way, we obtain the entities corresponding to image i_k .

3.2. Bidirectional Sub-Graph Sequentialization

With the list of entities in each image, we perform bidirectional sub-graph sequentialization to obtain a list of commonsense riddles in textual format, which can be readily used for contrastive training. The key idea is to find commonsense knowledge descriptions associated with each entity in the image, but with the entities’ names hidden by demonstrative pronouns. The pipeline of our graph operation can be summarized as follows:

1. To collect the commonsense knowledge associated with image i_k , we first query the directed knowledge graph to obtain sub-graphs where the nodes are connected, in either direction, to at least one entity in $\hat{\mathcal{N}}_k$.
2. Hide the names of the entities in the image by replacing the different subject nodes with “this” nodes.
3. We perform sequentialization to translate the sub-graph into a list of subject-hidden commonsense knowledge descriptions in textual format.

• **Bidirectional sub-graph query.** Specifically, we query the sub-graph that relates to i_k from the directed graph G , so that each edge within it is connected with at least a node representing one entity in $\hat{\mathcal{N}}_k$. The connection we need to check is bidirectional: both head and tail should be taken into account. Formally, we perform a bidirectional query to get sub-graphs,

$$\begin{aligned}
 G_k &= (V_k, E_k) \\
 \text{s.t. } V_k &= \left\{ v \mid v \in \ell_{\mathcal{E}}(u), u \in \hat{\mathcal{N}}_k \right\} \cup \hat{\mathcal{N}}_k, \\
 E_k &= \left\{ (u, r, w, v) \in \mathcal{E} \mid u \in \hat{\mathcal{N}}_k \right\} \\
 &\quad \cup \left\{ (v, r, w, u) \in \mathcal{E} \mid u \in \hat{\mathcal{N}}_k \right\},
 \end{aligned}$$

¹For example, NLTK [6] and SpaCy [21], two popular libraries for natural language processing in Python.

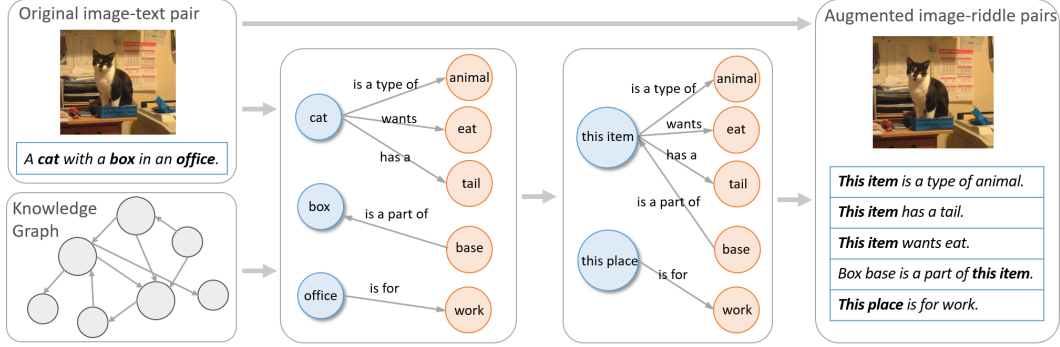


Figure 2. Illustration of VLCR construction process of DANCE.

and $\ell_{\mathcal{E}}(u)$ is the neighbors of node u when the edge set is \mathcal{E} . We end up with a sub-graph G_k . We refer the nodes that are directly from $\hat{\mathcal{N}}_k$ as subject nodes, i.e., $u \in \hat{\mathcal{N}}_k$.

• **Hiding subject names via node substitution.** After querying the sub-graphs, we perform node substitution that replaces subject nodes with “this” nodes, to hide the names of all the entities in the image. Specifically, we construct a mapping

$$f(\cdot): \mathcal{V} \rightarrow \mathcal{S},$$

that maps the actual entity nodes to a substitution set \mathcal{S} , which is defined as

$$\mathcal{S} = \{“this item”, “this person”, “this place”\}.$$

In detail, a node u is mapped to “this person” if it belongs to a “person” word, e.g., “lady”, “guy”, or mapped to “this place” if it has location property² or the subject name matches Places356 categories [89], or “this place” if neither of above is matched. Further, we filter out some edges in the sub-graph with weights below a certain threshold τ . More rigorously, the graph after substitution is

$$\begin{aligned} G'_k &= (V'_k, E'_k), \\ \text{s.t. } E'_k &= \left\{ (f(u), r, w, v) \mid (u, r, w, v) \in E_k, u \in \hat{\mathcal{N}}_k, w > \tau \right\} \\ &\quad \cup \left\{ (v, r, w, f(u)) \mid (v, r, w, u) \in E_k, u \in \hat{\mathcal{N}}_k, w > \tau \right\}, \\ V'_k &= \left\{ v \mid v \in \ell_{E'_k}(u), u \in \hat{\mathcal{N}}_k \right\} \cup \left\{ f(u) \mid u \in \hat{\mathcal{N}}_k \right\}. \end{aligned}$$

• **Riddle generation.** We concatenate the head, the relation, and the tail in each edge to generate riddles in natural language. Take an example from Fig. 2. For edge (“this item”, “is a type of”, 0.6, “animal”), we re-format it into natural language riddle “this item is a type of animal” and pair it back to image i_k . In this way, we obtain multiple image-riddle pairs for each image, which is readily usable for contrastive VL-model training. The whole process is performed automatically without human involvement. Leveraging mature image-text datasets and knowledge graphs, the quality of the generated data is well guaranteed.

²For example, in ConceptNet, such nodes has relation property “AtLocation”.

3.3. Training Strategy

Since image-riddle pairs are essentially image-text pairs, by mixing them with the existing VL databases with a certain ratio, we can pre-train or fine-tune the VL-model without changing the model architecture and other training settings. However, since the total amount of text generated by DANCE is several times larger than that of the existing VL dataset, simply merging our data with the original dataset will cause our data to dominate. Denote the proportion of the augmented data in the training batch as p . We observe that a larger p at the beginning and a smaller p in the later training stage can lead to good performance. Therefore, we adopt a curriculum learning strategy, with linearly decreasing p [65]. In this way, the percentage of original and our data sources can be controlled dynamically. There is no change to the inference stage.

4. Diagnostic Data and Automatic Evaluation

It is still an open problem to automatically and directly compare commonsense knowledge of VL-models without transferring to other downstream tasks. Thus, we introduce a diagnostic test set for comparison of the commonsense ability, in the form of a retrieval task compatible with various VL-models. Our task is divided into Text-Image and Image-Text retrieval. The former one is to retrieve which image best matches the description that requires commonsense but with the referred subject hidden, and the latter is vice versa. Thus, we mainly focus on the former one in the following paragraph of generation and evaluation. Formally, the model is either given a riddle d with a list of N_i candidate images $\mathcal{I} = \{i_1, \dots, i_{N_i}\}$ to choose from, or an image with a list of riddles. Models need to score the alignment between the images and the riddles and return a sorting of them. The data construction is based on the COCO test set and ConceptNet.

Generation of candidate set. Different from existing image-text datasets that usually contain one-to-one pairs, in the generated text set, there are multiple positive riddles for each image, and multiple positive images for each riddle. Thus, for test data, we also need to generate candidate sets

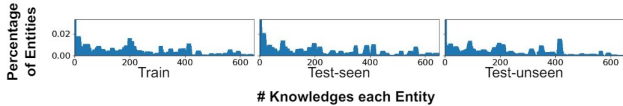


Figure 3. Visualization of the distributions of our train and Test-Image test splits.

including both positive and negative images with a consistent ratio. Suppose in the image list \mathcal{I} , there are n positive images i_1, \dots, i_n . They are chosen at random from the set of images that all contains the substituted entity in the riddle d . To ensure the high quality of the negative images, i.e., i_{n+1}, \dots, i_{N_i} , rather than random sampling from all possible negatives, we design to search for hard-negative samples. Specifically, We employ neighborhood hard-negative filtering, i.e., we find images whose entities are highly related to the subject entity, but none of these entities satisfy the riddle. To capture the correlation between entities, we use the graph distance in ConceptNet, i.e., from one entity, we filter for their nearest neighbor entities that are connected by the edges among three relationships: “*RelatedTo*”, “*Distinct-From*”, and “*Antonym*”. We construct the image-to-riddle data in the same way.

Generalization ability. To further diagnose the ability to infer new knowledge by using existing commonsense knowledge, we randomly hold out some knowledge from training. For example, given that “*pineapple can be found on a pizza*”, and “*pizza hut is one of the makers of pizza*”, we want to see whether the model can reason that “*pineapple may be required by pizza hut*”. Therefore, we further divide the test set into two splits: **test-seen** split, in which all knowledge appears in the training set, and **test-unseen** split where the corresponding relationships are not present in training, to see whether the model can reason about new knowledge using the existing ones. We also enforce that all the images in the two test splits are not present in training.

Automatic evaluation protocol. For automatic evaluation, we adopt perplexity score as the evaluation metric, following the works [13, 50]. In the experiment, we set the candidate number for each sample as 50, with the number of positive samples n between 1 to 15, and measure the average accuracy of retrieved ground truth inference, denoted as $Acc@50$. Our evaluation protocol is based on retrieval tasks, and therefore is compatible with most of the VL-model architectures.

Distribution of test-unseen split. We visualize the distribution of test-unseen split in comparison with training and test-seen split to verify that it represents a reasonable knowledge distribution. In Fig. 3, we show the distributions of the train, test-seen, and test-unseen based on COCO and ConceptNet. The x axis represents the number of common-sense knowledge descriptions associated with each entity, and the y axis represents the percentage of these entities among all entities. The distribution of the test-unseen split

Models		Text-Image		Image-Text	
		test-seen \uparrow	test-unseen \uparrow	test-seen \uparrow	test-unseen \uparrow
Random	-	0.2381	0.2380	0.2400	0.2362
Contr- astive	CLIP(ViT-L)	0.3951	0.3949	0.3817	0.3961
	OwIVit(ViT-L)	0.3673	0.3644	0.3325	0.3230
Matching	ViLT(ViLT-B)	0.4098	0.4077	0.3217	0.3534
	FLAVA(ViT-B)	0.4144	0.4093	0.3850	0.3843
Both	BLIP(ViT-L itm)	0.4030	0.4019	0.4017	0.4194
	BLIP(ViT-L itc)	0.3835	0.4007	0.3167	0.3100
	ALBEF(ViT-B)	0.3901	0.3792	0.3749	0.3832
Human	-	0.8202	0.8023	0.8497	0.8521

Table 1. Comparison with various state-of-the-art VL-models and human performance on our diagnostic test set.

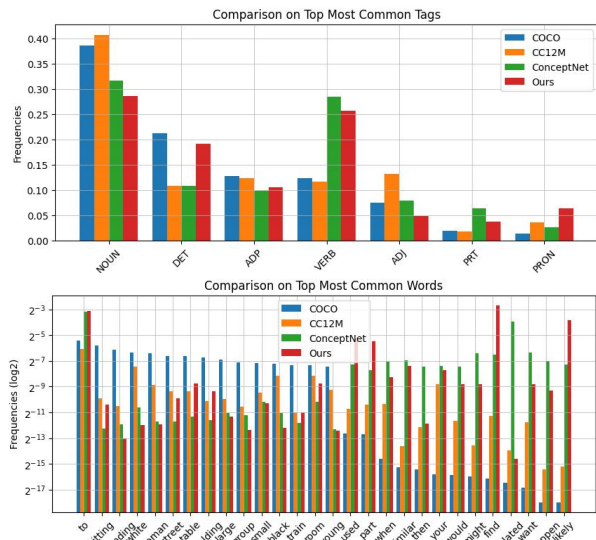


Figure 4. Comparison of the syntactic categories and words distributions of fundamental VL data (COCO [37] and CC12M [9]), NLP knowledge base (ConceptNet [66]) and ours. Commonsense is lacking in VL data, but it has significantly improved in ours.

does not shift much from the training and test-seen splits.

5. Experiment

In this section, we first highlight the commonsense lacking issue in both the popular VL datasets and the existing VL-models, then provide more analysis on the augmented training data, and finally provide detailed empirical evidence showing the effectiveness of the proposed DANCE strategy.

5.1. Commonsense Lacking Issue

Commonsense lacking in fundamental VL data. We find that the current fundamental VL datasets, on which VL-models are trained, are to blame for the commonsense knowledge-lacking issue. We show that the current fundamental datasets for VL-models can not provide sufficient commonsense knowledge compared to regular texts.

We illustrate the issue by comparing the most popular VL datasets (COCO [37] and CC 12M [9]) with the commonly used language-only knowledge bases (ConceptNet [66]) in terms of the distributions of the syntactic categories and words. In the upper part of Fig. 4, we compare the distributions of the most frequent part-of-speech (POS) tags with punctuation marks excluded. In the lower part, we show the comparison of the most frequent word tokens. There is a significant difference between top POS tag / word token distributions of VL datasets compared with those of the regular texts. We note that most frequent words in the text in existing VL datasets are nouns, especially for individual entities that appear in the images and the corresponding caption (e.g., “woman”, “train”, “room”). In contrast, the knowledge base ConceptNet has comparatively many more verbs, e.g., “used”, “would”, “find”, “want”, “happen”, “requires”, which contain richer information about the relationship between entities. In addition, the knowledge base includes more particles (PRT), like “to”, and pronouns (PRON) like “your”, which are associated with inter-connection information.

In order to develop common sense and reasoning ability, in addition to knowing each isolated entity, there is a high demand for rich but implicit information about the relationships and interconnections between entities. Thus, the fundamental VL dataset, which is primarily occupied by information about individual entities that appear explicitly, does not meet the requirements of VL-models for common knowledge, in terms of both learning or evaluation. This implies that we should enhance VL data with commonsense. Our augmented data via DANCE provides significantly more commonsense knowledge than the VL datasets. In Sec. 1 of the appendix, we include additional comparisons of VL data with common NLP data to further illustrate the commonsense lacking issue.

Baselines vs Human. Here we show the performance of various state-of-the-art VL-models on our diagnostic test set. Specifically, we consider VL-models in three categories: models trained by contrastive supervision, e.g., CLIP [52], OwlVit [44], by matching supervision like ViLT [27], FLAVA [62], and by the both, such as BLIP [33], ALBEF [34]. It is either impossible or difficult to directly test these models against knowledge-based benchmarks designed with other downstream tasks. As an additional reference for performance comparison, we also report the performance of random ordering as the lower bound, and human performance on a random sub-set with 50 samples each split as the upper bound. All the mentioned models are with their official checkpoints. CLIP is with ViT-L/14 backbone pre-trained on 400M images at 336-pixel resolution, ViLT is with ViLT-B/32 backbone pre-trained on 21M images and fine-tuned on COCO retrieval, and BLIP is with ViT-L/16 backbone pre-trained on 129M images with bootstrapping

	DANCE (this work)	VCR [86]	Visual- COMET [50]	OK-VQA [42]	S3VQA [22]
supervision	contrastive/matching	multi-choice	inference	VQA	VQA
# images	14.1M+	0.1M	59K	14K	7K
# texts	447M+	0.3M	1.5M	14K	7K
knowledge	general	people action	movie event	factoid	factoid

Table 2. Comparison with various knowledge-based datasets.

and find-tuned on COCO retrieval.

Results with our automatic evaluation metric $Acc@50$ are shown in Table 1. Through observation, we discover that, while most of the retrieval is easy for humans (83% on average, 81% for Text-Image, and 85% for Image-Text), they are hard for current state-of-the-art VL-models (<40% on average, 39.4% for Text-Image, and 36.3% for Image-Text) whose performances are only slightly better than the random result (24% on average).

5.2. Analysis on the Augmented Data

Implementation Details. Before the bidirectional query, we need to match the natural language words in text t_k to the knowledge graph entities. Therefore, we perform Unicode normalization, lower-casing, punctuation cleaning, underscoring and pre-pending to the words. For example, the English words “A traffic jam” becomes “/c/en/traffic.jam” by the standardization, so that it is matched to an entity. The threshold τ is set to 0.5. We are based on two mature human-annotated image-text paired datasets COCO [37], VG [28], and three web datasets SBU captions [49], Conceptual Captions (CC3M) [60], and Conceptual 12M (CC12M) [9], with 14M images in total. Our image splitting for COCO follows a popular split [3, 12, 20, 33, 54, 70, 73, 74, 80].

Dataset comparison with various knowledge-based datasets. In Table 2, we compare the training set generated by DANCE to relevant knowledge-based datasets and display their properties and statistics. We are the first commonsense knowledge-based dataset to focus on contrastive or matching supervision. We have larger-scale multi-source images and a corpus compared to the relevant datasets which are challenging to gather at scale, and we can expand even further if other image-text datasets or knowledge bases are included in the generation process. Also, our dataset includes various general knowledge from a consolidated commonsense knowledge graph, while the knowledge type in some of the relevant datasets (e.g., VCR, VisualCOMET) is limited to people or events in films.

Other statistics. The left part of Fig. 5 shows the distribution of commonsense knowledge type in the training set. We can observe that various types of commonsense knowledge are included in our generated data. The right part is the distribution of the average length of the commonsense riddles in our generated training set. The average riddle length generated from COCO and ConceptNet is 8.19.

Backbone	Pre-train	Fine-tune	Ours test set				COCO (5K test set)	
			Text-Image		Image-Text		TR@1↑	IR@1↑
			test-seen↑	test-unseen↑	test-seen↑	test-unseen↑		
ViT-B	14M	COCO	0.3986	0.3892	0.2916	0.3301	78.40	60.70
ViT-L	14M	COCO	0.4030	0.4019	0.4017	0.4194	81.12	63.96
ViT-B	14M+DANCE(part)	COCO	0.5107	0.5141	0.5252	0.5053	78.80	60.48
ViT-L	14M+DANCE(part)	COCO	0.5408	0.5326	0.5363	0.5113	80.89	64.15
ViT-L	14M+DANCE(whole)	COCO	0.5721	0.5458	0.5600	0.5242	81.92	65.26
ViT-B	14M	COCO+DANCE	0.4566	0.4077	0.3421	0.3845	77.94	60.83
ViT-L	14M	COCO+DANCE	0.4610	0.4333	0.4565	0.4395	81.86	64.17

Table 3. Effect of DANCE for pre-training (first five rows) and fine-tuning (last two rows), testing on ours test set and COCO retrieval.

Backbone	Pre-train	Fine-tune	OK-VQA Acc↑
ViT-B	14M	OK-VQA	29.45
ViT-B	14M+DANCE(part)	OK-VQA	37.56
ViT-L	14M	OK-VQA	33.14
ViT-L	14M+DANCE(part)	OK-VQA	38.55
ViT-L	14M+DANCE(whole)	OK-VQA	39.25

Table 4. Effect of DANCE for pre-training, testing on existing commonsense benchmark OK-VQA.

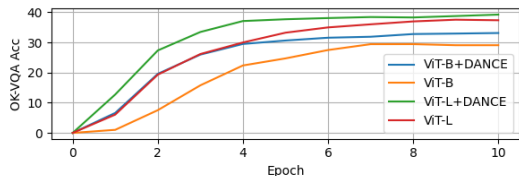


Figure 7. Performance on existing commonsense benchmark OK-VQA during fine-tuning.

domain, which means that there is no data leakage in the commonsense knowledge graph used for training.

The baselines are BLIP pre-trained on 14M images including COCO [37], VG [28], SBU captions [49], Conceptual Captions (CC3M) [60], and Conceptual 12M (CC12M) [9], and fine-tuned on OK-VQA. They are compared with DANCE pre-trained models respectively. Performing DANCE on a part of the pre-train data (COCO) generating 40M image-text pairs during the pre-training is denoted by DANCE(part), and on the full data generating 0.4B image-text pairs is denoted by DANCE(whole).

The results are shown in Table 4. Significant improvements can be seen by comparing the corresponding baselines and DANCE. When DANCE is applied just on part of pre-training data, the performances of ViT-Base and ViT-Large are increased by 3.69% and 1% respectively. And by using DANCE on the entire pre-train data, we achieve 1.69% improvement with ViT-Large. Besides, as shown in Fig. 7, we found that the model pre-trained on DANCE can achieve faster and more stable convergence than the ordinary pre-trained model on the OK-VQA benchmark.

DANCE for Fine-tuning. DANCE can also contribute positively in the fine-tuning stage. We evaluate the commonsense ability of the model on the proposed diagnostic test set, rather than on existing commonsense knowledge

DANCE Ratio	50%	30%	10%	0%	50→10%
OK-VQA Acc	30.32	32.48	32.03	29.45	33.14

Table 5. Ablation of proportions of DANCE-augmented data.

benchmarks, since the latter cannot evaluate models that are not fine-tuned on the specific downstream tasks, e.g., VQA. For fairness, DANCE and the corresponding baseline are fine-tuned on the same number of batches and steps, and on the same training set of COCO images. The results are reported in the last two rows of Table 3. Compared with the baseline, DANCE fine-tuning brings significant improvement in the commonsense test sets, and can still obtain comparable results on COCO.

Ablation. In the table 5, we study the impact of different proportions of DANCE-augmented data, i.e., p , on the performance in the downstream OK-VQA task. We find that a suitable ratio (30%) is beneficial to the performance improvement, while too high (50%) or too low (10%) reduces the performance improvement, and the curriculum learning strategy that linearly decreases the ratio from 50% to 10% achieves the best performance. The experiments are based on ViT-B backbone and 14M pre-train data.

6. Conclusion and Future Work

This paper takes a step towards injecting commonsense capability into VL-models. We first observed that VL-models are lacking commonsense ability as existing popular VL datasets do not contain much commonsense knowledge. Therefore, we propose a new training strategy DANCE which is compatible with most VL-models, by training on images paired with our generated entity-hidden commonsense riddles, in a scalable and automatic way. To support the commonsense evaluation of a suite of VL-models in a well-received way, a retrieval-based commonsense diagnostic benchmark is built. We then empirically verify the weaknesses of existing VL-models and the effectiveness of DANCE. Despite significant improvements in both the existing commonsense and our diagnostic benchmarks, we still face challenges. Towards human-like intelligence, awareness of commonsense knowledge is not enough. The model should be able to do reasoning, such as mathematical

and physical calculations in real-life scenarios. This is still weak in existing VL-models and is not included in existing commonsense knowledge bases. Future research could be conducted to analyze and improve various reasoning aspects of VL-models.

References

- [1] Conceptnet knowledge: Helium is used for filling party balloons. <https://conceptnet.io/c/en/helium?rel=/r/UsedFor&limit=1000>. Accessed: 2022-10-10. **7**
- [2] Junia Anacleto, Henry Lieberman, Marie Tsutsumi, Vânia Neris, Aparecido Carvalho, Jose Espinosa, Muriel Godoi, and Silvia Zem-Mascarenhas. Can common sense uncover cultural differences in computer applications? In *IFIP International Conference on Artificial Intelligence in Theory and Practice*, pages 1–10. Springer, 2006. **3**
- [3] 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*, pages 6077–6086, 2018. **6**
- [4] 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*, pages 2425–2433, 2015. **7**
- [5] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. Dbpedia: A nucleus for a web of open data. In *The semantic web*, pages 722–735. Springer, 2007. **2, 3**
- [6] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. ” O’Reilly Media, Inc.”, 2009. **3**
- [7] Min Cao, Shiping Li, Juntao Li, Liqiang Nie, and Min Zhang. Image-text retrieval: A survey on recent research and development. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5410–5417. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Survey Track. **2**
- [8] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. **2**
- [9] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021. **1, 2, 5, 6, 8, 13, 14**
- [10] Feilong Chen, Duzhen Zhang, Minglun Han, Xiuyi Chen, Jing Shi, Shuang Xu, and Bo Xu. Vlp: A survey on vision-language pre-training. *arXiv preprint arXiv:2202.09061*, 2022. **2**
- [11] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020. **2**
- [12] Wenliang Dai, Lu Hou, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. Enabling multimodal generation on clip via vision-language knowledge distillation. *arXiv preprint arXiv:2203.06386*, 2022. **6**
- [13] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 326–335, 2017. **5**
- [14] 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. **7, 13**
- [15] Xiaoyi Dong, Yinglin Zheng, Jianmin Bao, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, et al. Maskclip: Masked self-distillation advances contrastive language-image pretraining. *arXiv preprint arXiv:2208.12262*, 2022. **2**
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. **7**
- [17] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*, 2022. **3**
- [18] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2021. **13**
- [19] Wikimedia Foundation. Wikimedia downloads. **13, 14**
- [20] Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. Language models are general-purpose interfaces. *arXiv preprint arXiv:2206.06336*, 2022. **6**
- [21] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017. **3**
- [22] Aman Jain, Mayank Kothiyari, Vishwajeet Kumar, Preethi Jyothi, Ganesh Ramakrishnan, and Soumen Chakrabarti. Select, substitute, search: A new benchmark for knowledge-augmented visual question answering. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2491–2498, 2021. **2, 6**
- [23] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. **2**

- [24] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020. [2](#)
- [25] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetmodulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021. [2](#)
- [26] Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. I-bert: Integer-only bert quantization. In *International conference on machine learning*, pages 5506–5518. PMLR, 2021. [13](#)
- [27] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR, 2021. [1](#), [2](#), [6](#)
- [28] 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*, 123(1):32–73, 2017. [2](#), [6](#), [8](#)
- [29] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. *CoRR*, abs/1909.11942, 2019. [13](#)
- [30] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontañón. Fnet: Mixing tokens with fourier transforms. *CoRR*, abs/2105.03824, 2021. [13](#)
- [31] DB Lenat and RV Guha. Building large knowledge-based systems: Representation and inference in the cyc project. *Artificial Intelligence*, 61(1):4152, 1993. [3](#)
- [32] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11336–11344, 2020. [2](#)
- [33] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022. [1](#), [2](#), [6](#), [7](#)
- [34] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021. [2](#), [6](#), [7](#)
- [35] 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. [2](#)
- [36] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. *arXiv preprint arXiv:2110.05208*, 2021. [3](#)
- [37] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollr, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [1](#), [2](#), [5](#), [6](#), [8](#), [13](#), [14](#)
- [38] Yuanze Lin, Yujia Xie, Dongdong Chen, Yichong Xu, Chengguang Zhu, and Lu Yuan. Revive: Regional visual representation matters in knowledge-based visual question answering. *arXiv preprint arXiv:2206.01201*, 2022. [2](#)
- [39] 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. [13](#)
- [40] Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on real-life images with pre-trained vision-and-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2125–2134, 2021. [2](#)
- [41] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019. [2](#)
- [42] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [2](#), [6](#), [7](#)
- [43] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019. [2](#)
- [44] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022. [2](#), [6](#)
- [45] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. *arXiv preprint arXiv:2112.12750*, 2021. [3](#)
- [46] Sharan Narang, Hyung Won Chung, Yi Tay, Liam Fedus, Thibault Févry, Michael Matena, Karishma Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, et al. Do transformer modifications transfer across implementations and applications? In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5758–5773, 2021. [13](#)
- [47] David Amat Olóndriz, Pong Palau Puigdevall, and Adrià Salvador Palau. Foodi-ml: a large multi-language dataset of food, drinks and groceries images and descriptions. *arXiv preprint arXiv:2110.02035*, 2021. [2](#)
- [48] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. [2](#)

- [49] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011. 6, 8
- [50] Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. Visualcomet: Reasoning about the dynamic context of a still image. In *European Conference on Computer Vision*, pages 508–524. Springer, 2020. 3, 5, 6
- [51] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019. 2
- [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 1, 2, 6
- [53] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*, 2019. 13, 14
- [54] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018. 6
- [55] Negar Rostamzadeh, Seyedarian Hosseini, Thomas Boquet, Wojciech Stokowiec, Ying Zhang, Christian Jauvin, and Chris Pal. Fashion-gen: The generative fashion dataset and challenge. *arXiv preprint arXiv:1806.08317*, 2018. 2
- [56] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019. 13
- [57] Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Bideman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What language model to train if you have one million gpu hours? *arXiv preprint arXiv:2210.15424*, 2022. 13
- [58] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022. 2
- [59] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 2
- [60] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018. 2, 6, 8
- [61] Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Anna Rohrbach, Zhe Gan, Lijuan Wang, Lu Yuan, et al. K-lite: Learning transferable visual models with external knowledge. *arXiv preprint arXiv:2204.09222*, 2022. 2, 3
- [62] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. 2, 6
- [63] Push Singh et al. The public acquisition of commonsense knowledge. In *Proceedings of AAAI Spring Symposium: Acquiring (and Using) Linguistic (and World) Knowledge for Information Access*, 2002. 3
- [64] Amit Singhal et al. Introducing the knowledge graph: things, not strings. *Official google blog*, 5:16, 2012. 3
- [65] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. *International Journal of Computer Vision*, pages 1–40, 2022. 4
- [66] Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*, 2017. 2, 3, 5, 6, 13, 14
- [67] Robyn Speer and Catherine Havasi. Conceptnet 5: A large semantic network for relational knowledge. In *The People’s Web Meets NLP*, pages 161–176. Springer, 2013. 3
- [68] Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from pre-training and fine-tuning transformers. *arXiv preprint arXiv:2109.10686*, 2021. 13
- [69] Yi Tay, Vinh Q Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. Charformer: Fast character transformers via gradient-based subword tokenization. In *International Conference on Learning Representations*, 2021. 13
- [70] Yoad Tewel, Yoav Shalev, Idan Schwartz, and Lior Wolf. Zero-shot image-to-text generation for visual-semantic arithmetic. *arXiv preprint arXiv:2111.14447*, 2021. 6
- [71] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017. 2
- [72] Haoran Wang, Ying Zhang, Zhong Ji, Yanwei Pang, and Lin Ma. Consensus-aware visual-semantic embedding for image-text matching. 2020. 2
- [73] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luowei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. Omnivl: One foundation model for image-language and video-language tasks. *arXiv preprint arXiv:2209.07526*, 2022. 2, 6
- [74] Li Wang, Zechen Bai, Yonghua Zhang, and Hongtao Lu. Show, recall, and tell: Image captioning with recall mechanism. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 12176–12183, 2020. 6
- [75] Peng Wang, Qi Wu, Chunhua Shen, Anthony R Dick, and Anton van den Hengel. Explicit knowledge-based reasoning for visual question answering. In *IJCAI*, 2017. 2

- [76] Peng Wang, Qi Wu, Chunhua Shen, Anton Hengel, and Anthony Dick. Fvqa: Fact-based visual question answering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP, 06 2016. [2](#)
- [77] Hui Wu, Yupeng Gao, Xiaoxiao Guo, Ziad Al-Halah, Steven Rennie, Kristen Grauman, and Rogerio Feris. Fashion iq: A new dataset towards retrieving images by natural language feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11307–11317, 2021. [2](#)
- [78] Yichong Xu, Chenguang Zhu, Shuohang Wang, Siqi Sun, Hao Cheng, Xiaodong Liu, Jianfeng Gao, Pengcheng He, Michael Zeng, and Xuedong Huang. Human parity on commonsenseqa: Augmenting self-attention with external attention. *arXiv preprint arXiv:2112.03254*, 2021. [3](#)
- [79] Yichong Xu, Chenguang Zhu, Ruochen Xu, Yang Liu, Michael Zeng, and Xuedong Huang. Fusing context into knowledge graph for commonsense question answering. *arXiv preprint arXiv:2012.04808*, 2020. [3](#)
- [80] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10685–10694, 2019. [6](#)
- [81] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. *CoRR*, abs/1906.08237, 2019. [13](#)
- [82] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint arXiv:2111.07783*, 2021. [3](#)
- [83] Shuquan Ye, Dongdong Chen, Songfang Han, and Jing Liao. 3d question answering. *arXiv preprint arXiv:2112.08359*, 2021. [2](#)
- [84] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. [2](#)
- [85] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. [2](#)
- [86] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [3](#), [6](#)
- [87] Andy Zeng, Adrian Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari, Aavek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022. [1](#)
- [88] Si Zhang, Hanghang Tong, Jiejun Xu, and Ross Maciejewski. Graph convolutional networks: a comprehensive review. *Computational Social Networks*, 6(1):1–23, 2019. [2](#)
- [89] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. [4](#)
- [90] Chenguang Zhu, Yichong Xu, Xiang Ren, Bill Lin, Meng Jiang, and Wenhao Yu. Knowledge-augmented methods for natural language processing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 12–20, 2022. [3](#)

This appendix is organized as follows:

- In Section 7, we further illustrate the commonsense lacking issue by providing additional comparison of fundamental VL datasets with commonly used NLP data.
- In Section 8, we provide more visualizations of success examples of our method on the proposed diagnostic benchmark for both text-image and image-text retrieval.
- In Section 9, we provide more visualizations of success examples of our method on the OK-VQA benchmark.
- In Section 10, we summarize the statistics of the proposed diagnostic test data.
- In Section 11, we study the failure case of our DANCE augmented model.

7. Commonsense in Fundamental VL Data vs NLP Data

We further explore the commonsense lacking issue in the current fundamental VL data by comparing them with common natural language processing (NLP) data. Here we compare the distributions of the syntactic categories and words of the most popular VL datasets (COCO [37] and CC 12M [9]) with three commonly used NLP datasets: ConceptNet [66] the knowledge base dataset, Wikipedia [19] the popular [14, 29, 39, 56, 81] cleaned English-language articles with the size of 16GB, C4 [53] the popular used [18, 26, 30, 46, 57, 68, 69] English-language text sourced from the Common Crawl web scrape with the size of 745GB. The syntactic categories and word distributions comparison is shown in Fig. 8.

The upper part of Fig. 8 shows the distribution of the most frequent part-of-speech (POS) tags with punctuation marks excluded, and the lower part shows the most frequent word tokens. There is a significant difference between top POS tag/word token distributions of VL datasets compared with those of the regular texts. Similar to our observation in the main paper, the most frequent words in the text in existing VL datasets are nouns (NOUN) for **individual entities**, like “street”, “table”, “train”. In contrast, all the NLP datasets have apparently more verbs (VERB), like “have”, “used”, “find”, “want”, “happen” that contains richer information about the **relationship between entities**. Besides, the NLP datasets include more particles (PRT), like “to”, and pronouns (PRON) like “your”, which are associated with **interconnection** information. This further illustrates the lacking commonsense issue in the fundamental VL datasets.

While the implicit information about the **interconnections between entities** is in high demand for developing

commonsense and reasoning ability, the fundamental VL datasets are lacking it. This motivates us to use commonsense knowledge to improve VL data. In addition, the distribution of ours training data is also included for comparison. We can see that our data is similar to NLP data in terms of the interconnection between entities.

8. Additional Qualitative Results on Our Diagnostic Benchmark

In Fig. 9 and Fig. 10, we show additional qualitative comparison with the state-of-the-art VL-models on our diagnostic test set for text-image and image-text retrieval respectively. In Fig. 9, from left to right is the input text, the input images including a correct one (in blue) and two incorrect ones (in red), the scores by each individual model, and the commonsense knowledge from the knowledge graph [66] that required for retrieval. In Fig. 10, from left to right is the input image, the input texts including a correct one (in blue) and two incorrect ones (in red), the scores, and the related commonsense knowledge from the knowledge graph. We can see that all the baselines fail to identify the correct answers, which further illustrates the lacking of commonsense ability in the popular VL-models. In contrast, our DANCE pre-trained model successfully retrieves the correct ones. We note that all these images and the knowledge are held out from the training set. This further demonstrates the reasoning ability enhanced by our DANCE strategy.

9. Additional Qualitative Results on OK-VQA Benchmark

In Fig. 11, we show additional qualitative comparison with the state-of-the-art VL-models on the official validation split of the popular commonsense-aware OK-VQA dataset. We note that the validation split is not included during fine-tuning. From left to right is the input question, the input image, the answers by the baseline model BLIP, the DANCE pre-trained model and human, and the related commonsense knowledge from the knowledge graph. The baseline model struggles with these questions and predicts some relevant but wrong answers, which further demonstrates the lack of commonsense ability in the current VL-models. DANCE improves the VL-model’s commonsense ability in numerous aspects, including the commonsense knowledge of physics as shown in the first row, the commonsense of human behavior and motivation in the second and third rows, and the knowledge about animals in the fourth and fifth rows. This further demonstrates the commonsense ability enhanced by our DANCE strategy.

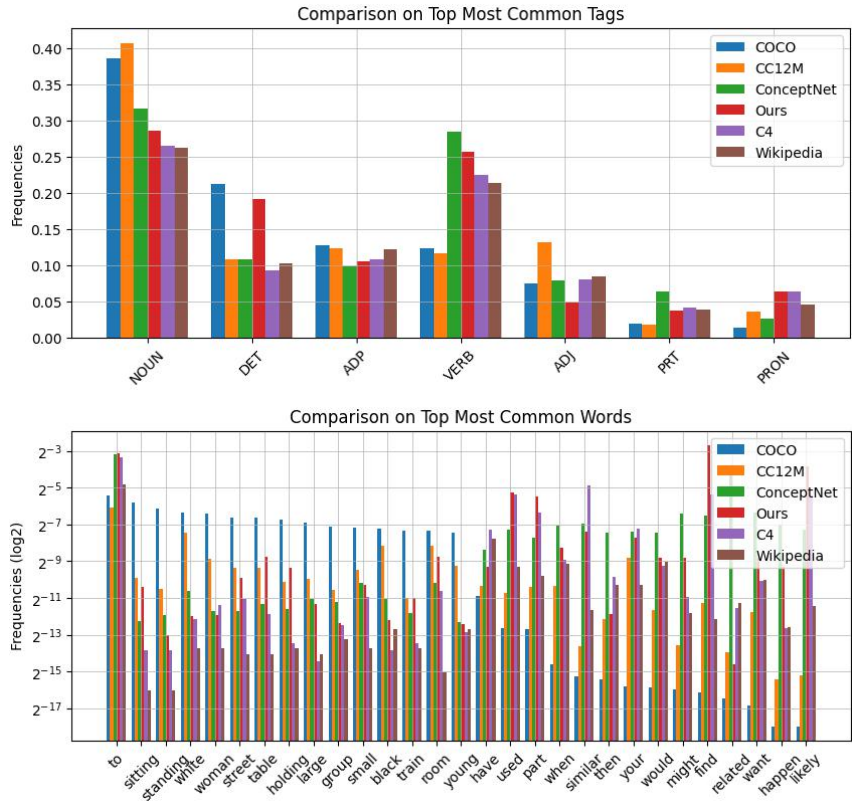


Figure 8. Comparison of the syntactic categories and words distributions of fundamental VL data (COCO [37] and CC12M [9]), ours training data generated by DANCE, and commonly used NLP data (ConceptNet [66], Wikipedia [19] and C4 [53]). Commonsense is lacking in VL data compared with NLP data, and is improved by DANCE strategy.



Figure 9. Qualitative examples from our diagnostic test set for text-image retrieval.



Figure 10. Qualitative examples from our diagnostic test set for image-text retrieval.

	Text-Image seen	Text-Image unseen	Image-Text seen	Image-Text unseen
# Images	4949	4974	500	500
# Texts	500	500	13930	14889
# Seen Images	0	0	0	0
# Seen Texts	500	0	13930	0

Table 6. Statistics of different splits of our diagnostic benchmark.

10. Statistics of Our Diagnostic Benchmark

In Table 6, we show the statistics of the four different splits of our diagnostic retrieval test set. Each row respectively represents the number of different images, the number of different text or riddles in each split, and the number of different images and texts that also appear in the training data. All these images for our test set does not appear in the training set. The knowledge in both Text-Image unseen split and Image-Text unseen split is held out from the training set.

11. Failure Case on OK-VQA Benchmark

In the main paper, we mainly focus on enhancing the VL-model’s ability to general commonsense via combining the VL data lacking commonsense with commonsense

knowledge graphs. However, our model learned from this commonsense-augmented data still suffers in some special real-life scenarios. Here we visualize the failure case of the model with DANCE pre-training in Fig. 12. The model fails to answer a question about counting or quantity. This indicates that the sense of numbers or the mathematical reasoning ability is still weak in existing VL-models, which is also not included in existing commonsense knowledge bases.


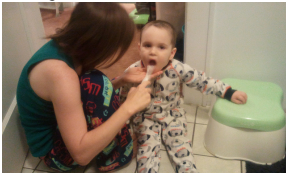

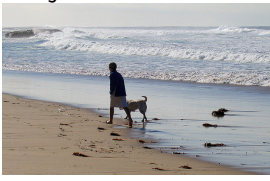

<p>Question: What celestial body controls the movements of the body of water featured in this photo?</p>	<p>Image: </p>	<p>✗ BLIP: wave ✓ Ours: moon Human: moon,moon,moon, moon,moon,moon, moon,moon,moon,moon</p>	<p>Commonsense knowledge: [[Ocean tides]] can be influenced by the [[moon]] [[The moon]] is for [[ocean tides]]</p>
<p>Question: This activity helps to ensure that what remains fresh?</p>	<p>Image: </p>	<p>✗ BLIP: toothpaste ✓ Ours: breath Human: brush teeth,brush teeth, brush teeth,brush teeth, breath,breath,breath, breath,breath,brush,brush</p>	<p>Commonsense knowledge: You will [[brush your teeth]] if you want to [[fresh your breath]]</p>
<p>Question: Should we go or stop?</p>	<p>Image: </p>	<p>✗ BLIP: slow down ✓ Ours: go Human: go,go,go,go,go, go,go,go,go,go</p>	<p>Commonsense knowledge: [[Green light]] means [[go ahead]]</p>
<p>Question: The animal in this image is said to be man's best what?</p>	<p>Image: </p>	<p>✗ BLIP: swim ✓ Ours: friend Human: friend,friend,friend, friend,friend,friend,best friend,best friend,dog,dog</p>	<p>Commonsense knowledge: [[A dog]] is [[a man's best friend]]</p>
<p>Question: Is this animal male or female?</p>	<p>Image: </p>	<p>✗ BLIP: male ✓ Ours: female Human: female,female,female, female,female,female,female, female,female,female</p>	<p>Commonsense knowledge: [[Rooster]] has [[a comb]]</p>

Figure 11. Qualitative examples from the commonsense-aware benchmark OK-VQA.

<p>Question: How many numbers are on this item?</p>	<p>Image: </p>	<p>✗ BLIP: 3 ✗ Ours: 24</p>	<p>Human: 12,12,12,12,12, 12,12,12,12,12</p>
---	---	---------------------------------	--

Figure 12. Case study of failure on the OK-VQA benchmark.