

Leveraging Visual Knowledge in Language Tasks: An Empirical Study on Intermediate Pre-training for Cross-modal Knowledge Transfer

Anonymous ACL submission

Abstract

Pre-trained language models are still far from human performance in tasks that need understanding of properties (e.g. appearance, measurable quantity) and affordances of everyday objects in the real world since the text lacks such information due to reporting bias. In this work, we study whether integrating visual knowledge into a language model can fill the gap. We investigate two types of knowledge transfer: (1) *text knowledge transfer* using image captions that may contain enriched visual knowledge and (2) *cross-modal knowledge transfer* using both images and captions with vision-language training objectives. On 5 downstream tasks that may need visual knowledge to solve the problem, we perform extensive empirical comparisons over the presented objectives. Our experiments show that visual knowledge transfer can improve performance in both low-resource and fully supervised settings.¹

1 Introduction

Pre-trained language models (PTLMs) such as BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019), and T5 (Raffel et al., 2020) have shown impressive results in various conventional natural language understanding (NLU) tasks by capturing syntactic and semantic knowledge from the pre-training tasks of *masked language modeling* and *masked span infilling* tasks on massive text corpora.

Though yielding good performance on various NLU downstream tasks, these pre-training objectives suffer from a lack of out-of-domain knowledge that is not explicitly present in the pre-training corpus (Gururangan et al., 2020a; Petroni et al., 2021; Schick and Schütze, 2020). Specifically, one type of knowledge that models often struggle with is the visual knowledge of common objects such as attributes (e.g. appearance, measurable quantity) and affordances. This is because this kind of knowledge is rarely explicitly described in the training

¹Code and data have been uploaded and will be published.

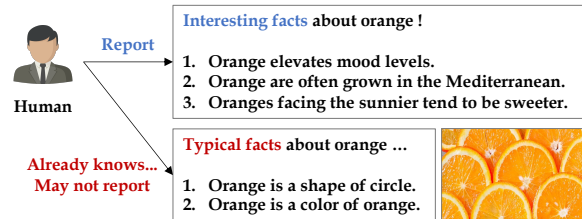


Figure 1: **Reporting Bias.** People tend to report what interests them rather than typical and general facts.

text due to reporting bias. For example, as shown in Figure 1, people tend to report what interests them rather than general facts such as a shape or color of oranges they already know.

Towards better knowledge-enhanced PTLMs, recent works incorporate external knowledge bases (e.g., knowledge graph, dictionary) to inject entity knowledge into PTLMs (Zhang et al., 2019; Peters et al., 2019; Wang et al., 2021; Yu et al., 2021) or retrieve knowledge from external knowledge bases to solve the problem (Lin et al., 2019; Wang et al., 2020). However, these approaches still suffer from a lack of visual knowledge that is important to understand the real world.

In this paper, we conduct systematic experiments to understand whether such visual knowledge can be transferred into LMs, and if so, how to perform effective knowledge transfer. Specifically, we look into a series of analysis question as follows: (1) Can intermediate pre-training (Pruksachatkun et al., 2020a) on image-caption pairs help transfer the knowledge? (2) What types of knowledge sources are more helpful? To answer questions, we explore various intermediate pre-training tasks (Pruksachatkun et al., 2020a) on two different sources: text-only (*text knowledge transfer* from visual domains) and image-caption pairs (*cross-modal knowledge transfer*).

For the text knowledge transfer, we utilize text corpus from visual domain, e.g., image captions. We leverage two training objectives for the lan-

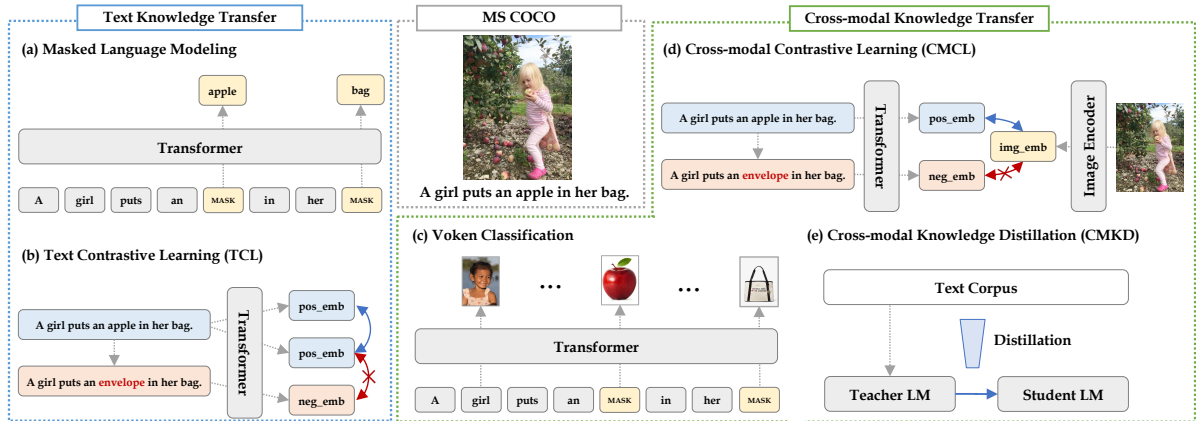


Figure 2: **Illustration of different methods for transferring visual knowledge into transformer-based language model.** In this example, we assume image-caption pair as an input. (a) *masked language model* (Devlin et al., 2018) on image captions. (b) *text contrastive learning* obtains positive example by dropout representation to learn better sentence representation while negative augmentation is optional. (c) *voken classification* employs token-level text-to-image retrieval to transfer visual knowledge. (d) *cross-modal contrastive learning* aims to train correct paring of images and captions. (e) *cross-modal knowledge distillation* transfers knowledge from the teacher model, which is trained by cross-modal contrastive learning, into student model.

072 guage model: (1) *masked language modeling* fol- 103
 073 lows the domain adaptive pre-training scheme (Gu- 104
 074 rurangan et al., 2020a), assuming the corpus con- 105
 075 tains enriched visual knowledge or physical com- 106
 076 monsense knowledge; (2) *text contrastive learning* 107
 077 augments the sentence representation with dropout 108
 078 to create positive samples while considering all 109
 079 others in the batch as negative samples for the con- 110
 080 trastive learning (Gao et al., 2021), assuming train- 111
 081 ing better sentence representations leads to better 112
 082 understanding of the corpus. 113

083 For the cross-modal knowledge transfer, we ex- 114
 084 plore multiple methods to transfer visual-related 115
 085 knowledge to LMs: (1) *masked language model- 116
 086 ing with visual clues* incorporates visual clues to 117
 087 capture dependencies between visual and linguis- 118
 088 tic contents (Su et al., 2019); (2) *voken classifica- 119
 089 tion* contextually aligns language tokens to their 120
 090 related images (called "vokens") to transfer visual 121
 091 knowledge into LMs (Tan and Bansal, 2020); (3) 122
 092 *cross-modal contrastive learning* aims to improve 123
 093 text representations by maximizing the agreement 124
 094 between correct image-text pairs versus random (in- 125
 095 batch) and adversarial negative pairs by contrastive 126
 096 learning between image and text modalities; and 127
 097 (4) *cross-modal knowledge distillation* transfers 128
 098 the knowledge from the teacher model, which is 129
 099 trained by cross-modal contrastive learning on im- 130
 100 age and text modalities, to the student language 131
 101 model using knowledge distillation.

102 We perform comprehensive comparisons on

103 five downstream tasks that may require visual 104
 105 or physical commonsense knowledge, including 106
 107 PIQA (Bisk et al., 2020), Visual Paraphrasing 108
 109 (VP) (Lin and Parikh, 2015), CSQA (Talmor et al., 110
 111 2018), OBQA (Mihaylov et al., 2018), and Rid- 112
 113 dleSense (Lin et al., 2021). Results suggest that: 114
 115 (1) Simple intermediate pre-training on captions 116
 117 can help improving performance on commonsense 118
 119 reasoning that needs physical or visual knowl- 120
 121 edge. (2) Cross-modal knowledge transfer approaches 122
 123 consistently improve the performance in a large 124
 125 margin when only few train examples are available. 126
 127 (3) Cross-modal contrastive learning shows that it 128
 129 is best for packaging visual knowledge into LMs. 130
 131

2 Analysis Setup 117

118 In this work, we study how to transfer the visual 119
 120 knowledge into language models. For this study, 121
 122 we introduce our analysis setup: problem formula- 123
 124 tion, analysis questions, and knowledge corpora. 125
 126

2.1 Problem Formulation 122

123 We focus on a pre-trained text encoder f_L and 124
 125 an image encoder f_V if images are available. f_L 126
 127 and f_V are initialized with pre-trained model and 128
 129 we continue to pre-train the models on different 130
 131 sources and tasks, which we call *intermediate pre- 132
 133 training* (Gururangan et al., 2020b; Pruksachatkun 134
 135 et al., 2020b). After the intermediate pre-training, 136
 137 we fine-tune f_L on downstream NLU tasks. Ex- 138
 139 isting NLU benchmarks have been trained against 140
 141

standard supervised learning paradigms that typically require a large number of question answering examples which need a large annotation efforts. However, in scenarios where the number of labeled examples is small, the model tends to overfit the training examples and shows poor generalization performance on test set. Here, we evaluate the intermediate pre-training objective’s generalization ability on test set in both fully supervised and low-resource settings.

2.2 Analysis Questions

In this paper, we provide a comprehensive study for transferring the visual knowledge into LMs. Visual knowledge transfer can be done in two approaches, depending on the source to be trained: (1) *Text knowledge transfer* using the text corpus in the visual domain, e.g., image captions and (2) *cross-modal knowledge transfer* which passes visual knowledge about common objects to LMs by training over paired image and captions. By evaluating the model on 5 downstream datasets that require physical and visual commonsense knowledge, we explore following three research questions.

Q1: Can intermediate pre-training on external knowledge sources help transfer visual knowledge to augment text encoders? We investigate diverse intermediate pre-training methods with external knowledge sources including caption data to inject visual information from images and captions into LMs. We first analyze the performance of text and cross-modal knowledge transfer methods with a image-caption dataset, and we additionally study text knowledge transfer methods with other text corpora such as GenericsKB (Bhakhavatsalam et al., 2020), Wiki103 (Merity et al., 2016) and BookCorpus (Zhu et al., 2015a).

Q2: What types of knowledge sources are more helpful for visual knowledge transfer? As mentioned above, we have two categories to exploit visual information: (1) *text knowledge transfer* and (2) *cross-modal knowledge transfer*. Here, we explore which type of knowledge transfer is more useful to transfer the visual knowledge into LMs.

Q3: What intermediate pre-training objectives are effective for cross-modal knowledge transfer? We present three pre-training objectives for cross-modal knowledge transfer: (1) token classification, (2) contrastive learning, and (3) knowledge distillation. Here, we want to present which strategy is best suited for cross-modal knowledge

Dataset	# Train	# Dev	# Test	# choices
PIQA	14,113	1,838	2,000	2
VP	21,988	2,000	6,057	2
CSQA	8,500	1,221	1,241	5
OBQA	4,957	500	500	4
RiddleSense	3,510	1,021	1,202	5

Table 1: **Downstream task data statistics.** We create in-house test set for PIQA and CSQA, and in-house dev set for VP by splitting the train set.

transfer. Furthermore, we study how to enhance cross-modal contrastive learning with adversarial negative samplings.

2.3 Pre-training Data

To transfer the visual knowledge, we collect 250K image-caption pairs from MS COCO (Lin et al., 2014; Chen et al., 2015). MS COCO contains images reflecting the composition of actual everyday scenes and corresponding captions which describe contextual reasoning between objects in the scene. We only use captions for text knowledge transfer while we use both images and captions for cross-modal knowledge transfer. As an ablation study, we explore other text corpora such as GenericsKB (Bhakhavatsalam et al., 2020), Wiki103 (Merity et al., 2016) and BookCorpus (Zhu et al., 2015a).

2.4 Downstream Tasks and Datasets

For downstream benchmarks, we find tasks that can benefit from visual knowledge: multiple choice question answering tasks including PIQA (Bisk et al., 2020) which requires physical commonsense reasoning, CSQA (Talmor et al., 2018) for general understanding of commonsense reasoning, OBQA (Mihaylov et al., 2018) that needs elementary-level science knowledge, and RiddleSense (RS) (Lin et al., 2021) for complex understanding of figurative language, and binary classification task including Visual Paraphrasing (VP) (Lin and Parikh, 2015) that needs scene understanding. We use in-house test sets made from training sets for PIQA and CSQA since test set is not provided to public. We list the data statistics in Table 1. Moreover, We additionally test on GLUE (Wang et al., 2018) to evaluate the general text understanding.

2.5 Evaluation Protocol

We evaluate the models in both fully supervised and low-resource settings. For both settings, we

consider accuracy for 5 different classification tasks and get average performance over tasks to check the final performance. In the fully supervised setting, we evaluate models with 3 different random seeds and report the average accuracy. In the low-resource setting, we set the size of the train data to 64 or 128. For each experiment, we run over 5 different sub-samples and show the average accuracy.

3 Method

In this section, we introduce the following two approaches to integrate visual knowledge into LMs: (1) *text knowledge transfer*; and (2) *cross-modal knowledge transfer*. Throughout this section, we assume the data is a collection of image x^v and caption x^l pairs $\{(x_i^v, x_i^l)\}_{i=1}^m$ (m is the size of the pairs) and image encoder f_V and text encoder f_L are given. Note that we use the same text encoder.

3.1 Text Knowledge Transfer

For text knowledge transfer, we investigate following pre-training objectives: (1) *masked language modeling*; and (2) *text contrastive learning*.

Masked Language Modeling (MLM) Following BERT (Devlin et al., 2018), we select 15% of input tokens and replace them with [MASK]. Of the selected tokens, 80% are replaced, 10% are not changed and 10% are replaced by random vocabulary token. Here, we employ dynamic masking, which performs random masking and replacement during training to prevent the same masking for the same examples (Liu et al., 2019). MLM objective is the cross-entropy loss for masked token predictions :

$$\ell_{\text{MLM}}(x_i^l) = -\log p(x_i^l | x^{\text{masked}}), \quad (1)$$

where x_i is the i -th token and x^{masked} is a mask.

Text Contrastive Learning (TCL) Contrastive learning aims to learn representations by pulling positive pairs closer and pushing negative pairs apart. Here, we employ the contrastive framework with cross-entropy objective and in-batch negatives (Chen et al., 2020a; Gao et al., 2021). Given a text encoder f_L , and a caption x_i^l , we first get text representations using the encoders $h_i^l = f_L(x_i^l)$. Following Gao et al. (2021), we create identical positive sample h_i^{l+} by different dropout representations. The contrastive loss is defined as follows:

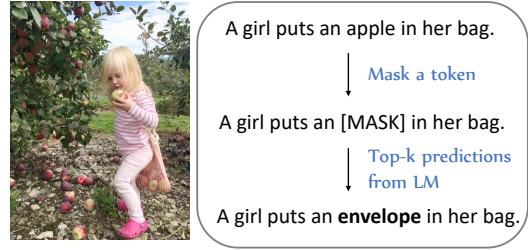


Figure 3: **LM perturbation.** We create adversarial negatives using language models.

$$\ell_i^l = -\log \frac{e^{\text{sim}(h_i^l, h_i^{l+})/\tau}}{\sum_{j=1}^N e^{\text{sim}(h_i^l, h_j^l)/\tau}}, \quad (2)$$

where N is a batch size and $\text{sim}(\cdot)$ represents cosine similarity, i.e., $\text{sim}(u, v) = u \cdot v / \|u\| \|v\|$. τ represents a temperature parameter.

3.2 Cross-modal Knowledge Transfer

Language models might learn additional information from visual sources such as images and captions. So we include a variety of vision-based approaches and investigate the approaches whether they can benefit from visual sources. We introduce vision-based approaches as follows.

Voken Classification Vokenization (Tan and Bansal, 2020) employs token-level text-to-image retrieval to transfer visual knowledge. It aligns language tokens to their related images (called “vokens”) to transfer visual knowledge into LMs, and call it “voken classification”. Given text x and a voken v_i for the i -th token, the loss is defined as

$$\ell_i^{\text{voken}} = -\log(p(v_i | x)). \quad (3)$$

Similar to masked language modeling, it classifies each token to a corresponding voken. Vokenization trains language models with the voken classification task and MLM.

Masked Language Modeling with Visual Clues VL-BERT (Su et al., 2019) adopts masked language modeling with visual clues in which models are given a caption with masked tokens and an image and predict the masked tokens using visual clues. VL-BERT is pre-trained on Conceptual Captions (Sharma et al., 2018) as an image-caption corpus, and BooksCorpus (Zhu et al., 2015b) and English Wikipedia as text-only corpora. It shows its effectiveness in many vision-language tasks. We investigate whether this model also succeed in NLP tasks and compare it with others.

Cross-modal Contrastive Learning (CMCL)

To harness the visual knowledge from image-caption datasets, we adopt contrastive loss on image and text vectors. Given an image encoder f_V , a text encoder f_L , and an image-caption pair (x_i^v, x_i^l) , we first get image and text representations using the encoders $h_i^v = f_V(x_i^v)$, $h_i^l = f_L(x_i^l)$. Then the contrastive learning objective contains two loss functions: an image-to-text contrastive loss $\ell^{(v,l)}$ and a text-to-image contrastive loss $\ell^{(l,v)}$. The image-to-text contrastive loss is defined as follows:

$$\ell_i^{(v,l)} = -\log \frac{e^{\text{sim}(h_i^v, h_i^l)/\tau}}{\sum_{j=1}^N e^{\text{sim}(h_i^v, h_j^l)/\tau}}, \quad (4)$$

where N is a batch size and $\text{sim}(\cdot)$ represents cosine similarity. This loss encourages a closer distance between representations of aligned image-caption pairs than unaligned pairs given an image and multiple captions. Similarly, the text-to-image contrastive loss $\ell^{(l,v)}$ is defined as follows:

$$\ell_i^{(l,v)} = -\log \frac{e^{\text{sim}(h_i^l, h_i^v)/\tau}}{\sum_{j=1}^N e^{\text{sim}(h_i^l, h_j^v)/\tau}}. \quad (5)$$

The final loss is defined as

$$L = \frac{1}{N} \sum_{i=1}^N (\ell_i^{(v,l)} + \ell_i^{(l,v)}). \quad (6)$$

CLIP (Radford et al., 2021) and ConVIRT (Zhang et al., 2020) also adopt contrastive learning, but we freeze the image encoder in training and use the trained text encoder for downstream tasks.

CMCL with Adversarial Negative Samples (ANS)

As in-batch negatives in CMCL are not challenging enough for models to distinguish, we present adversarial negative sampling strategy to improve CMCL. Given an image-caption pair (x_i^v, x_i^l) , we define a LM-perturbed sentence x_i^{l-} , which is a hard negative where n is replaced with a different word n' from a probability distribution of PTLMs. We expect the l^- is syntactically correct and plausible sentence even the word n is replaced to n' , while it does not semantically match to the corresponding image x_i^v . With such hard negative, we try to make more challenging task so that models can effectively learn from the task. For example, we choose a word ‘girl’ in the sentence ‘A girl puts an apple in her bag.’ in Figure 3. Then we mask the word with [MASK] token to do masked token predictions by PTLMs. Then we get top- k predictions from language models and replace

the masked tokens with one of the predicted ones. To avoid false negative sentences which may have the same semantics as the original sentence, we introduce an additional filtering step: if the masked predictions are synonyms or hypernyms of the original tokens, we discard the predictions. We use WordNet (Miller, 1995) to find synonyms and hypernyms. The contrastive loss with hard negative is defined as follows:

$$-\log \frac{e^{\text{sim}(h_i^v, h_i^l)/\tau}}{\sum_{j=1}^N e^{\text{sim}(h_i^v, h_j^l)/\tau} + \sum_{k=1}^M e^{\text{sim}(h_i^v, h_k^{l-})/\tau}}, \quad (7)$$

where M is the number of hard negative samples per positive pair. This formula is only for image-to-text contrastive loss $\ell^{(v,l)}$ and final loss is defined to same as equation (6).

CMCL with Positive Sample Augmentation (PSA)

In ANS, we filter perturbed sentences where the masked predictions are synonyms or hypernyms of the original tokens. Instead of excluding these perturbed sentences, another option is to include them as additional positive samples l^+ to the paired images. We name this as positive sample augmentation (PSA). It also adopts LM-perturbed negative samples as in ANS.

Cross-modal Knowledge Distillation (CMKD)

Cross-modal knowledge distillation is to transfer knowledge between different modalities, e.g., image modality and text modality. In this category, CMKD is to transfer knowledge from a teacher model which is knowledgeable about visual information. VidLanKD (Tang et al., 2021) also utilizes a cross-modal knowledge distillation method to help with general language understanding. A teacher model is first trained using contrastive learning on a video-text dataset, and then it transfers its knowledge to a student language model using KD on a text corpus. Their contrastive learning loss (hinge loss) is defined as

$$L = \sum_i^N [\max(0, \alpha - \text{sim}(h_i^v, h_i^l) + \text{sim}(h_i^{v'}, h_i^l)) + \max(0, \alpha - \text{sim}(h_i^v, h_i^{l'}) + \text{sim}(h_i^v, h_i^{l'}))], \quad (8)$$

where v' and l' are a random image and caption text, respectively. α is the margin between the similarities of a positive pair and a negative pair. Instead of video datasets, we use a MS COCO dataset to train a teacher model and use two versions of contrastive learning, equations (6) and (8).

	Model	PIQA		VP		CSQA		OBQA		RiddleSense		Average	
		64	128	64	128	64	128	64	128	64	128	64	128
-	BERT-base	52.6±0.9	53.8±0.1	85.9±1.1	86.6±0.7	35.8±0.7	37.8±0.3	31.3±1.2	32.0±0.7	24.7±0.1	25.2±0.2	46.1	47.1
Caption	MLM	53.1±0.2	54.3±0.3	86.5±0.3	87.3±0.4	35.7±0.3	36.7±0.1	33.4±0.6	34.2±0.3	26.3±0.1	26.5±0.2	47.0	47.8
	TCL	52.6±0.5	52.9±0.6	86.4±0.1	88.0±0.1	35.7±0.2	36.1±0.3	34.2±1.4	35.2±0.7	30.3±0.5	30.7±0.4	47.8	48.5
	TCL + MLM	53.6±0.7	54.6±0.2	84.2±0.2	87.6±0.3	33.6±2.2	35.1±0.6	31.8±2.3	34.3±0.5	20.6±0.0	20.6±0.0	44.7	46.4
	TCL + ANS	50.0±0.7	50.5±0.6	67.3±0.4	68.2±0.7	26.8±1.2	27.5±0.5	33.4±1.1	35.0±1.0	26.1±1.7	26.5±1.8	40.7	41.5
	TCL + PSA + ANS	51.1±0.1	51.2±0.4	66.0±0.0	66.0±0.0	22.7±0.9	22.9±0.1	30.2±3.1	31.8±0.4	23.5±1.2	25.2±1.5	38.7	39.4
Caption-Image Pairs	VL-BERT-base	53.1±0.6	53.9±0.4	88.5±0.3	88.4±0.5	36.2±0.7	36.8±0.8	33.4±1.2	34.6±1.2	26.1±0.8	26.1±0.9	47.7	48.5
	Vokenization	50.5±0.5	51.1±0.4	68.8±1.6	78.1±1.9	19.2±1.4	21.5±0.8	31.2±2.7	33.2±2.2	17.1±0.5	16.7±0.7	37.3	40.1
	VidLanKD	55.0±0.4	55.6±0.5	86.7±0.5	88.5±0.5	37.1±1.0	38.6±0.5	31.8±1.3	32.6±1.0	24.4±0.0	24.4±0.0	47.0	47.9
	VidLanKD variant	55.3±0.3	55.2±0.4	87.4±0.1	88.2±0.6	37.3±1.2	38.9±0.5	32.4±2.1	32.2±1.1	24.4±0.0	24.4±0.0	47.3	47.7
	CMKD (VL-BERT-large)	54.7±0.5	54.5±0.2	86.5±0.8	88.4±0.4	36.7±0.4	38.5±0.4	29.8±0.8	31.7±0.2	25.2±0.1	25.2±0.0	46.5	47.6
	CMCL	54.7±0.4	55.1±0.1	87.9±0.3	88.9±0.2	36.3±0.3	38.4±0.4	31.1±1.1	32.8±0.9	25.0±0.2	25.4±0.4	47.0	48.1
	CMCL + ANS	55.4±0.1	55.7±0.2	88.1±0.9	88.9±0.7	37.5±0.8	39.0±0.2	32.2±0.7	32.0±0.6	27.4±0.0	27.5±0.1	48.1	48.6
	CMCL + PSA + ANS	55.4±0.2	55.1±0.2	88.8±1.0	88.2±0.2	37.0±0.3	38.1±0.3	34.1±0.4	34.8±0.9	26.7±0.4	28.8±0.7	48.4	49.0

Table 2: **Performance (accuracy) in low-resource setting.** We test models on diverse datasets with low-resource learning (64 and 128 training samples). We use captions in the MS COCO dataset for text knowledge transfer methods and images and captions for cross-modal knowledge transfer methods. We get average performance on 64 and 128 training samples. **Bold** and underlined numbers refer to the best and second-best performance, respectively.

As another version of CMKD, we consider distilling visual knowledge from a pre-trained vision-language model, VL-BERT, which is knowledgeable about grounded language. We adopt masked language modeling on Wikitext103 (Merity et al., 2016), a subset of English Wikipedia, in the knowledge distillation step. For knowledge distillation, we adopt Neuron Selectivity Transfer (NST) (Huang and Wang, 2017), which proves the effectiveness in VidLanKD (Tang et al., 2021).

4 Experimental Settings

For all the approaches, we use bert-base-uncased (Devlin et al., 2018) as text encoder f_L and ResNeXt101 (Xie et al., 2017) as an image encoder f_V . We continue to pre-train the encoders in our experiments. For text knowledge transfer, (1) MLM follows the exact setting of codebase in huggingface² which uses dynamic masking strategy to conduct language modeling task. (2) TCL conducts contrastive learning with f_L . We choose the best checkpoint by the best spearman correlation on STSb (Cer et al., 2017). For cross-modal knowledge transfer, (1) CMKD explores VL-BERT, Vokenization, and VidLanKD approaches. Here, we use VL-BERT-large model to do CMKD. We use the VL-BERT and Vokenization checkpoints from their official codebases³. VidLanKD trains a teacher model by two versions of contrastive learning (equations (6) and (8)) on MS COCO dataset. We set $\alpha = 1$ in VidLanKD (equation (8)). (2) CMCL conducts

²<https://github.com/huggingface/transformers/tree/master/examples/pytorch/language-modeling>

³<https://github.com/jackroos/VL-BERT>, <https://github.com/airsplay/vokenization>

	Model	PIQA	VP	CSQA	OBQA	RiddleSense	Average
-	BERT-base	62.5±1.3	93.1±0.4	53.2±1.2	52.2±0.5	38.9±0.9	59.9
Caption	MLM	63.8±0.9	93.5±0.1	52.6±0.3	53.9±1.1	39.3±1.4	60.6
	TCL	62.1±0.5	93.5±0.4	49.0±0.5	54.1±1.0	41.2±0.3	60.1
	TCL + MLM	62.3±0.7	93.2±0.3	49.0±0.4	49.0±0.8	40.5±0.5	58.8
	TCL + ANS	60.1±1.2	93.3±0.1	47.0±0.1	50.2±0.9	36.7±0.8	57.4
	TCL + PSA + ANS	59.5±1.0	92.4±0.3	34.0±1.3	44.6±1.4	28.4±2.3	51.7
Caption-Image Pairs	VL-BERT-base	63.8±1.5	93.6±0.1	50.3±1.1	49.6±2.3	39.1±1.0	59.2
	Vokenization	58.4±5.1	92.7±0.3	45.0±0.2	48.1±0.8	33.5±0.7	55.5
	VidLanKD	63.1±1.1	93.7±0.4	52.4±0.8	50.6±3.9	39.5±1.7	59.8
	VidLanKD variant	64.1±0.2	93.8±0.3	53.6±0.5	47.9±4.3	38.8±2.0	59.6
	CMKD (VL-BERT-large)	63.8±0.0	93.7±0.7	53.3±1.4	48.7±3.0	38.7±0.4	59.6
	CMCL	62.7±0.1	93.3±0.3	50.8±0.9	52.3±0.7	37.6±1.0	59.2
	CMCL + ANS	63.5±0.1	93.3±0.3	50.3±0.1	52.9±0.3	38.4±0.9	59.7
	CMCL + PSA + ANS	63.9±0.5	94.3±0.1	50.9±0.3	52.4±1.2	39.0±0.3	60.1

Table 3: **Performance (accuracy) in fully supervised setting.** **Bold** and underlined numbers refer to the best and second-best performance, respectively.

contrastive learning with f_L and f_V . Here, we set $\tau = 0.05$ (equations (4) and (5)). (3) CMCL with ANS chooses three noun words or verb words to do masked prediction and use top-5 predictions from f_L as replacement. We filter out synonyms and hypernyms of original words using WordNet (Miller, 1995). (4) CMCL with PSA includes the perturbed sentences with synonyms and hypernyms as additional positive samples. In CMCL, we adopt ResNeXt101 (Xie et al., 2017) as an image encoder f_V and BERT as a text encoder f_L . TCL and CMCL train with batch size 64, maximum sequence length 20, learning rate $1e-4$ for 3 epochs. For fine-tuning on downstream tasks, we do grid search on learning rates $\{5e-5, 1e-4, 3e-4, 4e-4, 5e-4, 6e-4\}$ and choose the best learning rate. We set maximum epochs to 30 in low-resource and 15 in fully supervised settings.

5 Results and Analysis

We analyze the main results of intermediate pre-training. Tables 2 and 3 show the main results of

	Model	RTE	MRPC	STS-B	CoLA	SST-2	QNLI	QQP	Avg.
	- BERT-base	70.0	<u>87.9</u>	89.1	57.4	91.3	90.4	89.3	82.3
Caption	MLM	62.8	87.0	89.1	53.9	92.6	91.1	90.9	81.0
	TCL	58.4	83.1	88.2	55.5	<u>91.9</u>	91.4	90.9	79.9
	TCL + MLM	54.8	81.6	87.2	53.6	<u>91.9</u>	90.9	89.2	78.5
	TCL + ANS	56.3	83.9	87.0	51.5	91.3	<u>91.2</u>	<u>89.4</u>	78.6
	TCL + PSA + ANS	52.3	75.6	81.5	17.4	90.0	85.8	88.2	70.1
Caption-Image Pairs	VL-BERT-base	57.4	85.7	<u>89.5</u>	58.1	90.6	89.7	88.7	80.0
	Vokenization	53.0	87.0	83.3	51.3	91.4	89.2	88.5	77.7
	VidLanKD	67.5	87.8	89.4	<u>57.7</u>	90.7	90.3	88.6	<u>81.7</u>
	VidLanKD variant	68.5	87.9	89.7	54.9	91.1	90.5	88.6	81.6
	CMKD (VL-BERT-large)	68.5	88.5	89.3	55.4	90.9	89.7	88.6	81.6
	CMCL	63.5	82.5	89.5	51.1	90.4	90.0	88.4	79.3
	CMCL + ANS	69.6	86.8	89.4	56.1	90.7	90.5	88.6	<u>81.7</u>
	CMCL + PSA + ANS	<u>69.8</u>	86.2	89.0	55.3	90.4	90.5	88.6	81.6

Table 4: **Performance (accuracy) on GLUE benchmark.** Bold and underlined numbers refer to the best and second-best performance, respectively.

low-resource learning and fully supervised learning with the MS COCO captioning dataset, respectively. We train the models with a few training examples, 64 and 128, to understand the better initialization. We argue that if a model obtains better performance in the low-resource setup, then it is a faster learner and has better generalization on downstream tasks.

Can text intermediate pre-training help improve text encoders? Text intermediate pre-training using MLM and TCL on a caption corpus improves the performance on downstream tasks in both low-resource and fully supervised settings. In particular, TCL shows significant improvement on OBQA and RiddleSense over BERT (p-value < 0.01). These results suggest that text intermediate pre-training on visual-related datasets helps performance on commonsense reasoning tasks.

Can cross-modal intermediate pre-training help transfer visual knowledge to augment text encoders? We observe that cross-modal intermediate pre-training is helpful in both fully supervised and low-resource settings (See Table 2 and 3). Specifically, CMKD with VidLanKD variant outperforms the baseline by 1.6% point on the PIQA dataset in fully supervised setting. CMCL also shows its effectiveness. However, we could find that it becomes more powerful when equipped with PSA and ANS. It suggests that data augmentation for positive and negative sampling is an important factor for CMCL. In low-resource setting, we find that cross-modal knowledge transfer helps better initialization and lets models learn new tasks faster.

What intermediate pre-training objectives are effective for cross-modal knowledge transfer? Among various cross-modal knowledge transfer methods, we study which method is the most effective for cross-modal knowledge transfer. Overall,

CMCL with PSA and ANS shows the best performance among all cross-modal methods. Interestingly, VL-BERT also shows better performance than BERT-base on all datasets in the low-resource setting. This suggests that exploiting images in masked language modeling task help transfer the knowledge to language models.

What types of knowledge sources are most helpful? Here, we investigate whether using an image source in addition to a text source can further improve the model. To answer this question, we analyze methods from different types of sources: text-only and text-image pair sources. We focus on the methods that use the contrastive learning objective: TCL and CMCL. Note that these two methods share the same objective but CMCL trains on cross modalities which are images and captions while TCL only trains on captions. Overall, TCL performs slightly better than CMCL in low-resource and fully supervised settings. Interestingly, additional negative samples (ANS) and positive samples in TCL decreases the performance while they help CMCL to improve the performance. We conjecture that perturbed sentences in ANS might not be semantically negative to the original sentence so models learn from wrong labels.

5.1 Ablation Study

How do models perform on general NLU tasks? Table 4 presents results on GLUE benchmark. In GLUE, text intermediate pre-training methods slightly underperform the original BERT-base. We conjecture that the intermediate pre-training on caption data might sacrifice knowledge of general language understanding.

Analysis on diverse text corpora Table 5 represents text approaches with different pre-training corpora: MS COCO captions (Lin et al., 2014; Chen et al., 2015), GenericsKB (Bhaktavatsalam et al., 2020), BooksCorpus (Zhu et al., 2015a), and WikiText103 (Merity et al., 2016). We sample 250k sentences from each corpus for a fair comparison. We notice that caption datasets are useful on OBQA and RiddleSense datasets while GenericsKB are the most helpful on PIQA datasets. Results are expected since GenericsKB contains a lot of everyday statements that contain various types of commonsense.

Different training sizes. We test different training sizes on PIQA in Fig. 4. In the experiment,

Model	PIQA			VP			CSQA			OBQA			RiddleSense		
	64	128	Full	64	128	Full	64	128	Full	64	128	Full	64	128	Full
- BERT-base	52.6 \pm 0.9	53.8 \pm 0.1	62.5 \pm 1.3	85.9 \pm 1.1	86.6 \pm 0.7	93.1 \pm 0.4	35.8 \pm 0.7	37.8 \pm 0.3	53.2 \pm 1.2	31.3 \pm 1.2	32.0 \pm 0.7	52.2 \pm 0.5	24.7 \pm 0.1	25.2 \pm 0.2	38.9 \pm 0.9
CP: MLM	53.1 \pm 0.2	54.3 \pm 0.3	63.8 \pm 0.9	86.5 \pm 0.3	87.3 \pm 0.4	93.5 \pm 0.1	35.7 \pm 0.3	37.7 \pm 0.1	52.6 \pm 0.3	33.4 \pm 0.6	34.2 \pm 0.3	53.9 \pm 1.1	26.3 \pm 0.1	26.5 \pm 0.2	39.3 \pm 1.4
TCL	52.6 \pm 0.5	52.9 \pm 0.6	62.1 \pm 0.5	86.4 \pm 0.1	88.0 \pm 0.1	93.5 \pm 0.4	35.7 \pm 0.2	36.1 \pm 0.3	49.0 \pm 0.5	34.2 \pm 1.4	35.2 \pm 0.7	54.1 \pm 1.0	30.3 \pm 0.5	30.7 \pm 0.4	41.2 \pm 0.3
GK: MLM	53.2 \pm 0.1	53.6 \pm 0.4	64.9 \pm 0.1	86.2 \pm 0.9	87.6 \pm 0.3	93.0 \pm 0.3	34.6 \pm 0.7	35.3 \pm 1.3	51.6 \pm 0.5	31.7 \pm 0.9	32.3 \pm 1.0	53.1 \pm 0.9	25.8 \pm 0.6	26.3 \pm 0.1	39.3 \pm 0.7
TCL	56.0 \pm 1.0	56.4 \pm 0.2	64.4 \pm 0.1	88.9 \pm 0.7	89.4 \pm 0.2	93.3 \pm 0.5	37.8 \pm 0.5	38.7 \pm 0.5	51.0 \pm 0.5	31.7 \pm 0.9	32.3 \pm 1.0	52.6 \pm 0.8	27.4 \pm 0.2	28.1 \pm 0.7	40.9 \pm 0.8
BC: MLM	54.1 \pm 0.3	54.1 \pm 0.8	63.3 \pm 0.6	86.4 \pm 0.8	87.5 \pm 0.5	93.0 \pm 0.3	29.8 \pm 0.8	32.1 \pm 0.9	50.8 \pm 0.3	29.6 \pm 0.8	31.4 \pm 0.7	50.2 \pm 0.4	22.6 \pm 0.0	22.7 \pm 0.0	36.7 \pm 1.3
TCL	52.4 \pm 0.1	53.1 \pm 0.4	63.1 \pm 0.3	87.1 \pm 1.9	89.7 \pm 0.1	93.2 \pm 0.2	38.0 \pm 0.5	38.1 \pm 1.1	51.5 \pm 0.1	33.8 \pm 2.7	34.0 \pm 2.1	55.6 \pm 0.4	28.9 \pm 0.4	29.1 \pm 0.3	41.2 \pm 2.3
WT: MLM	52.7 \pm 0.2	53.0 \pm 0.3	63.8 \pm 0.6	85.3 \pm 2.8	88.1 \pm 0.3	93.5 \pm 0.1	33.2 \pm 1.4	34.6 \pm 0.5	52.5 \pm 0.2	32.4 \pm 2.3	33.0 \pm 0.7	52.3 \pm 0.3	24.4 \pm 0.0	24.4 \pm 0.0	39.4 \pm 2.0
TCL	52.9 \pm 0.9	53.4 \pm 0.4	62.7 \pm 0.6	67.3 \pm 0.6	68.6 \pm 0.7	93.3 \pm 0.3	31.3 \pm 1.6	32.4 \pm 0.7	48.2 \pm 0.3	31.5 \pm 3.5	33.1 \pm 0.6	53.0 \pm 0.0	24.8 \pm 1.3	24.8 \pm 0.6	36.3 \pm 1.0

Table 5: **Results of text knowledge transfer methods with different corpora.** We pre-train text knowledge transfer methods, MLM and TCL, with different corpora. CP is MS COCO captions, GK is GenericsKB, BC is BooksCorpus, and WT is WikiText. **Bold** and underlined numbers refer to the best and second-best performance, respectively.

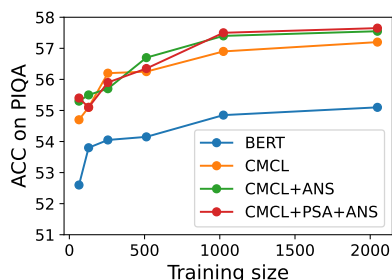


Figure 4: **Results on varying training sizes.** We test methods with different training sizes.

we observe that CMCL consistently outperforms BERT on all training sizes. Additional negative sample (ANS) improves the CMCL on different training sizes, and positive sample augmentation boosts the performance of CMCL further. This suggests including perturbed sentences as positive and negative samples are useful to cross-modal knowledge transfer.

6 Related Work

Text Knowledge enhanced methods. Recently, huge efforts on integrating knowledge into PTLMs have been made. One typical form of knowledge is a knowledge graph. There have been efforts of using knowledge graph to inject entity and relation representations, which are pre-computed from external source, into PTLMs (Zhang et al., 2019; Peters et al., 2019; He et al., 2020; Phang et al., 2020). Some other works try to retrieve or generate the sub-graph from the graph to solve the problem (Lin et al., 2019; Wang et al., 2020). Another existing form of knowledge is extra large-scale corpus. Works that use such corpus present knowledge-related pre-training objectives such as concept order recovering (Zhou et al., 2021), entity category prediction (Yu et al., 2020) and source of knowledge prediction (Wang et al., 2021; Calixto et al., 2021). They are mostly focused on inject-

ing world knowledge presented in text, rather than physical and visual commonsense knowledge that can be found in images.

Cross-modal knowledge enhanced methods.

There is a extensive line of works for a variety of vision-language tasks, such as VL-BERT (Su et al., 2019), VisualBert (Li et al., 2019), and Uniter (Chen et al., 2020b). These models aim to improve vision-language tasks, e.g., VQA (Goyal et al., 2017), and they are found to be not effective in improving language tasks (Tan and Bansal, 2020). Another line of works is to transfer visual knowledge to language models: Vokenization (Tan and Bansal, 2020) and VidLanKD (Tang et al., 2021). Vokenization employs token-level text-to-image retrieval to transfer visual knowledge to language models. For this, Vokenization introduces 30k vokens and matches each token into the limited voken space; it may have approximation errors. VidLanKD adopts contrastive learning to train a teacher model on video datasets and uses distillation approaches to distill visual knowledge from the teacher to a student model.

7 Conclusion

We study whether intermediate pre-training on visual knowledge can help transfer visual knowledge into LMs. We investigate text knowledge transfer and cross-modal knowledge transfer using images and captions. In our empirical analysis, we observe that intermediate pre-training on captions can help improving performance and cross-modal knowledge transfer approaches consistently improve performance. When the transfer methods are equipped with additional positive and negative samples, they show better performance. Future works include improving both commonsense reasoning and general language understanding.

References

- 595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
- Sumithra Bhakthavatsalam, Chloe Anastasiades, and Peter E. Clark. 2020. Genericskb: A knowledge base of generic statements. *ArXiv*, abs/2005.00660.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7432–7439.
- Iacer Calixto, Alessandro Raganato, and Tommaso Pasini. 2021. Wikipedia entities as rendezvous across languages: Grounding multilingual language models by predicting wikipedia hyperlinks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3651–3661.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. **SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation**. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020a. **A simple framework for contrastive learning of visual representations**. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020b. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer.
- 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*.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6904–6913.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020a. **Don’t stop pretraining: Adapt language models to domains and tasks**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020b. **Don’t stop pretraining: adapt language models to domains and tasks**. *arXiv preprint arXiv:2004.10964*.
- Bin He, Di Zhou, Jinghui Xiao, Xin Jiang, Qun Liu, Nicholas Jing Yuan, and Tong Xu. 2020. **BERT-MK: Integrating graph contextualized knowledge into pre-trained language models**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2281–2290, Online. Association for Computational Linguistics.
- Zehao Huang and Naiyan Wang. 2017. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. **KagNet: Knowledge-aware graph networks for commonsense reasoning**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2829–2839, Hong Kong, China. Association for Computational Linguistics.
- Bill Yuchen Lin, Ziyi Wu, Yichi Yang, Dong-Ho Lee, and Xiang Ren. 2021. Riddlesense: Reasoning about riddle questions featuring linguistic creativity and commonsense knowledge. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1504–1515.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.
- Xiao Lin and Devi Parikh. 2015. Don’t just listen, use your imagination: Leveraging visual common sense for non-visual tasks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2984–2993.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- 650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706

707	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. <i>arXiv preprint arXiv:1809.02789</i> .	
708		
709		
710		
711	George A Miller. 1995. Wordnet: a lexical database for english. <i>Communications of the ACM</i> , 38(11):39–41.	
712		
713	Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. 2019. Knowledge enhanced contextual word representations. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 43–54, Hong Kong, China. Association for Computational Linguistics.	
714		
715		
716		
717		
718		
719		
720		
721		
722	Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. KILT: a benchmark for knowledge intensive language tasks. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2523–2544, Online. Association for Computational Linguistics.	
723		
724		
725		
726		
727		
728		
729		
730		
731		
732	Jason Phang, Iacer Calixto, Phu Mon Htut, Yada Pruksachatkun, Haokun Liu, Clara Vania, Katharina Kann, and Samuel R Bowman. 2020. English intermediate-task training improves zero-shot cross-lingual transfer too. <i>arXiv preprint arXiv:2005.13013</i> .	
733		
734		
735		
736		
737	Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel R. Bowman. 2020a. Intermediate-task transfer learning with pretrained language models: When and why does it work? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5231–5247, Online. Association for Computational Linguistics.	
738		
739		
740		
741		
742		
743		
744		
745		
746	Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel R Bowman. 2020b. Intermediate-task transfer learning with pretrained models for natural language understanding: When and why does it work? <i>arXiv preprint arXiv:2005.00628</i> .	
747		
748		
749		
750		
751		
752		
753	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. <i>arXiv preprint arXiv:2103.00020</i> .	
754		
755		
756		
757		
758		
759	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
760		
761		
762		
763		
764		
	Timo Schick and Hinrich Schütze. 2020. Rare words: A major problem for contextualized embeddings and how to fix it by attentive mimicking. In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 8766–8774. AAAI Press.	765 766 767 768 769 770 771 772 773 774
	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2556–2565.	775 776 777 778 779 780
	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In <i>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17</i> , page 4444–4451. AAAI Press.	781 782 783 784 785
	Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. Vi-bert: Pre-training of generic visual-linguistic representations. <i>arXiv preprint arXiv:1908.08530</i> .	786 787 788 789
	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. Commonsenseqa: A question answering challenge targeting commonsense knowledge. <i>arXiv preprint arXiv:1811.00937</i> .	790 791 792 793
	Hao Tan and Mohit Bansal. 2020. Vokenization: improving language understanding with contextualized, visual-grounded supervision. <i>arXiv preprint arXiv:2010.06775</i> .	794 795 796 797
	Zineng Tang, Jaemin Cho, Hao Tan, and Mohit Bansal. 2021. Vidlankd: Improving language understanding via video-distilled knowledge transfer. <i>arXiv preprint arXiv:2107.02681</i> .	798 799 800 801
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	802 803 804 805 806 807 808 809
	Peifeng Wang, Nanyun Peng, Filip Ilievski, Pedro Szekely, and Xiang Ren. 2020. Connecting the dots: A knowledgeable path generator for commonsense question answering. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4129–4140, Online. Association for Computational Linguistics.	810 811 812 813 814 815 816
	Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. 2021. K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters.	817 818 819 820

821 In *Findings of the Association for Computational*
822 *Linguistics: ACL-IJCNLP 2021*, pages 1405–1418,
823 Online. Association for Computational Linguistics.

824 Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu,
825 and Kaiming He. 2017. Aggregated residual transfor-
826 mations for deep neural networks. In *Proceedings of*
827 *the IEEE conference on computer vision and pattern*
828 *recognition*, pages 1492–1500.

829 Donghan Yu, Chenguang Zhu, Yiming Yang, and
830 Michael Zeng. 2020. Jacket: Joint pre-training of
831 knowledge graph and language understanding. *arXiv*
832 *preprint arXiv:2010.00796*.

833 Wenhao Yu, Chenguang Zhu, Yuwei Fang, Donghan Yu,
834 Shuohang Wang, Yichong Xu, Michael Zeng, and
835 Meng Jiang. 2021. Dict-bert: Enhancing language
836 model pre-training with dictionary. *arXiv preprint*
837 *arXiv:2110.06490*.

838 Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christo-
839 pher D Manning, and Curtis P Langlotz. 2020.
840 Contrastive learning of medical visual representa-
841 tions from paired images and text. *arXiv preprint*
842 *arXiv:2010.00747*.

843 Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang,
844 Maosong Sun, and Qun Liu. 2019. [ERNIE: En-](#)
845 [hanced language representation with informative en-](#)
846 [tities](#). In *Proceedings of the 57th Annual Meeting of*
847 *the Association for Computational Linguistics*, pages
848 1441–1451, Florence, Italy. Association for Compu-
849 tational Linguistics.

850 Wangchunshu Zhou, Dong-Ho Lee, Ravi Kiran Sel-
851 vam, Seyeon Lee, and Xiang Ren. 2021. [Pre-training](#)
852 [text-to-text transformers for concept-centric common](#)
853 [sense](#). In *International Conference on Learning Rep-*
854 *resentations*.

855 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhut-
856 dinov, Raquel Urtasun, Antonio Torralba, and Sanja
857 Fidler. 2015a. Aligning books and movies: Towards
858 story-like visual explanations by watching movies
859 and reading books. In *The IEEE International Con-*
860 *ference on Computer Vision (ICCV)*.

861 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhut-
862 dinov, Raquel Urtasun, Antonio Torralba, and Sanja
863 Fidler. 2015b. Aligning books and movies: Towards
864 story-like visual explanations by watching movies
865 and reading books. In *Proceedings of the IEEE in-*
866 *ternational conference on computer vision*, pages
867 19–27.

868 **A Dataset Properties**

869 PIQA is a multiple-choice question answering task,
870 which chooses the most appropriate solution for
871 physical commonsense questions, which may need
872 illustration or description of physical interaction in
873 the real world. VP is to tell if two descriptions are
874 describing the same scene or two different scenes.
875 While they seem like purely textual tasks, they re-
876 quire visual common sense to answer. CSQA is
877 a multiple-choice question answering task that re-
878 quires commonsense reasoning to answer. It is built
879 from ConceptNet ([Speer et al., 2017](#)). OBQA is
880 a multiple-choice question answering task, which
881 is modeled after open book exams on elementary-
882 level core science questions. The task generally
883 requires open book fact but also additional com-
884 monsense which can be learnt from scientific illus-
885 tration. RiddleSense is a multiple-choice riddle-
886 style question answering which requires complex
887 commonsense reasoning ability and understanding
888 of figurative language which may benefit from vi-
889 sual knowledge.