

A Multi-Modal Context Reasoning Approach for Conditional Inference on Joint Textual and Visual Clues

Yunxin Li^{1*}, Baotian Hu^{1†}, Xinyu Chen¹, Yuxin Ding¹, Lin Ma², Min Zhang¹

¹Harbin Institute of Technology, Shenzhen, China, ²Meituan, Beijing

{hubaotian, yxding, zhangmin2021}@hit.edu.cn

{liyunxin987, chenxinyuhitsz}@163.com, forest.linma@gmail.com

Abstract

Conditional inference on joint textual and visual clues is a multi-modal reasoning task that textual clues provide prior permutation or external knowledge, which are complementary with visual content and pivotal to deducing the correct option. Previous methods utilizing pretrained vision-language models (VLMs) have achieved impressive performances, yet they show a lack of multimodal context reasoning capability, especially for text-modal information. To address this issue, we propose a **Multi-modal Context Reasoning** approach, named *ModCR*. Compared to VLMs performing reasoning via cross modal semantic alignment, it regards the given textual abstract semantic and objective image information as the pre-context information and embeds them into the language model to perform context reasoning. Different from recent vision-aided language models used in natural language processing, *ModCR* incorporates the multi-view semantic alignment information between language and vision by introducing the learnable alignment prefix between image and text in the pretrained language model. This makes the language model well-suitable for such multi-modal reasoning scenario on joint textual and visual clues. We conduct extensive experiments on two corresponding data sets and experimental results show significantly improved performance (exact gain by 4.8% on PMR test set) compared to previous strong baselines. Code Link: <https://github.com/YunxinLi/Multimodal-Context-Reasoning>.

1 Introduction

Cross modal reasoning is a hot research topic both in natural language processing and computer vision communities. Most cross modal reasoning tasks, such as Visual Question Answering (Antol et al., 2015; Wu et al., 2017; Shah et al., 2019;

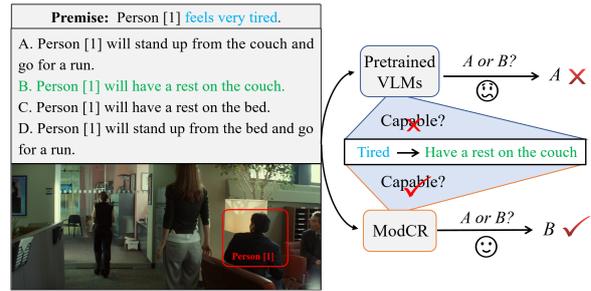


Figure 1: A case from the PMR (Dong et al., 2022) data set, where the correct option is answer **B**. The blue-color words represent the pivotal textual clue to infer the correctness of answers A and B.

Yusuf et al., 2022), Visual Dialog (Zhang et al., 2022; Chen et al., 2022), Visual Entailment, (Xie et al., 2019; Do et al., 2020) and Visual Commonsense Reasoning (Zellers et al., 2019a; Ye and Kovashka, 2021; Li et al., 2022a), concentrate on the visual reasoning scenario that relies primarily on image information. The given text (or question) is highly attached to the image and lacks prior permutation, e.g., the common question “Why is person 4 pointing to person 1” shown in VCR (Zellers et al., 2019a) data set. For another practical cross modal reasoning scenario (Dong et al., 2022), the textual modality often provides prior permutation or complementary information with the source image, such as the commonsense knowledge, and the personalities, feelings, or relationships of persons, as the premise shown in Figure 1. In this paper, we focus on such conditional inference on joint textual and visual clues, where the specific task form is to select the correct option from the candidate set according to the given textual premise and image.

Previous methods (Chen et al., 2020; Krojer et al., 2022; Li et al., 2020; Dong et al., 2022; Wang et al., 2022) usually input the concatenated sequence of textual premise, image, and candidate answer into powerful pretrained vision-language models (VLMs) and employ a task-specific classi-

*†Corresponding author.

fier to infer the result with attention to the joint representation obtained from VLMs. Although these methods work well for reasoning based mainly on visual clues, they suffer from one major shortcoming: the reasoning process does not fully utilize the abstract semantic information of given premise text to perform in-context reasoning. As the case shown in Figure 1, pretrained VLMs know “*person [1] sits on the couch, not the bed*” from the image, yet struggle to effectively infer that the person will “*have a rest on the couch*” according to “*feels very tired*” presented in the premise. It may be attributed to that pretrained VLMs mostly map different modalities into a unified space (Long et al., 2022) and perform cross modal semantic alignment and fusion. They neglect the in-context learning based on the given multi-modal semantics of language and vision during pertaining, like next sentence prediction. Fortunately, pretrained language models (PLMs) such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), BART (Lewis et al., 2020), and GPT3 (Brown et al., 2020), are powerfully capable of in-context learning and have achieved successful performance on natural language inference and open-ended text generation. PLMs can infer the next-step intent according to the given abstract text information compared to pretrained VLMs. Hence, we propose a simple and effective Multi-modal In-Context Reasoning approach named *ModCR* for this multi-modal reasoning task, taking advantages of VLMs and PLMs.

Specifically, *ModCR* employs a pretrained visual encoder equipped with a vision mapping network to obtain the image representation and convert it into the learnable visual prefix. The visual prefix and textual premise are regarded as two types of pre-context. They will be fed to the in-context reasoner, i.e., language model, to infer the correctness of answer. Considering the semantic gap between visual prefix and text in the language model, we first utilize a multi-grained vision-language semantic aligner to gain the multi-view alignment representation between image and text. Afterwards, we devise an alignment mapping network to capture the pivotal alignment information and convert it into the learnable cross-modal alignment prefix. Finally, we fed the two prefixes, premise, and answer into the language model to perform cross modal reasoning in the instruction template-based slot-filling method. In this way, *ModCR* bridges the semantic gap between visual content and text in

the language model through introducing the cross-modal alignment prefix. It makes use of the abstract semantic of premise and objective image information via the self-attention mechanism in PLMs.

To verify the effectiveness of *ModCR*, we conduct extensive experiments on two cross modal reasoning data sets: PMR (Dong et al., 2022) and VCR (Zellers et al., 2019a). The experimental results show that the proposed method significantly outperforms previous strong baselines. The ablation and case studies indicate that *ModCR* is capable of in-context reasoning based on multi-modal information.

Our contributions can be summarised as follows:

- We propose a multi-modal in-context reasoning framework for conditional inference on joint textual and visual clues, utilizing the in-context learning capability of PLMs.
- To the best of our knowledge, we are the first to introduce the multi-view alignment information between vision and language into the language model to perform cross modal reasoning, bridging the semantic gap between vision and language in PLMs.
- Experimental results show that *ModCR* achieves state-of-the-art performance on two corresponding data sets. It significantly outperforms previous vision-aided language models and pretrained VLMs-based approaches.

2 Related Works

Pretrained VLMs for Cross Modal Reasoning.

Cross modal reasoning (Chen et al., 2021; Long et al., 2022) is a challenging task that requires a cross modal understanding of images and texts with relational reasoning to infer the correct option. Vision-language models are thus proposed to represent, align, and fuse the image and text information and perform task-specific reasoning such as Visual Question Answering (Antol et al., 2015; Wu et al., 2017; Shah et al., 2019; Yusuf et al., 2022; Gao et al., 2022), Visual Dialog (Zhang et al., 2022; Chen et al., 2022; Lin and Byrne, 2022) or Storytelling (Huang et al., 2016; Yu et al., 2021b), Visual Entailment, (Xie et al., 2019; Do et al., 2020), Visual Commonsense Reasoning (Zellers et al., 2019a; Ye and Kovashka, 2021; Li et al., 2022a). Over the past few years, significant performance has been made for developing vision-language

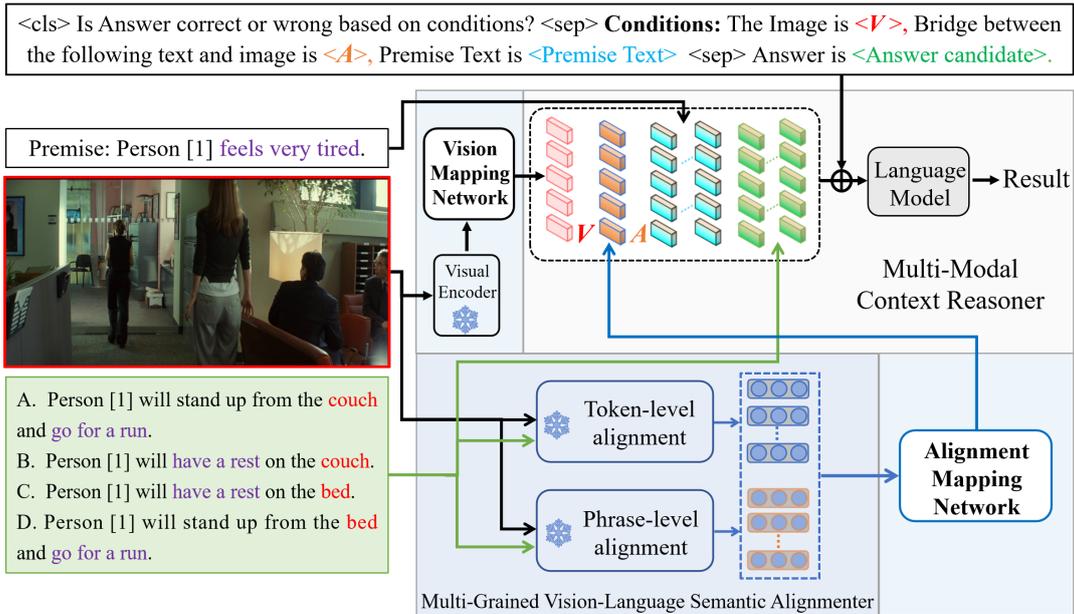


Figure 2: The overall workflow of ModCR. The top part presents the slot-filling instruction template used in the multi-modal in-context reasoner. The purple words show the relevant content between the premise and answers. The red words in answers are related to the image information. “V” and “A” indicate the vectors of visual and cross-modal alignment prefixes, respectively.

models, owing to the Transformer (Vaswani et al., 2017) architecture and large-scale multi-modal web data (Bugliarello et al., 2021; Lin et al., 2021). These pretrained VLMs could be divided into single-stream (Wang et al., 2021; Li et al., 2021) and double-stream (Radford et al., 2021; Jia et al., 2021; Lu et al., 2022a) types according to multi-modal information interaction methods. Our work explores how to expand and ameliorate pretrained VLMs to conditional inference on joint textual and visual clues.

Vision-aided Language Models. Images can provide explicit and diverse visual information to improve the imaginative representation of language. Recent works show that vision-aided language models have achieved promising performance on natural language understanding (Lu et al., 2022b) and open-ended text generation tasks (Zhu et al., 2022) such as text completion (Zellers et al., 2019b), story generation (Fan et al., 2018), and concept-to-text (Barzilay and Lapata, 2005). Some works (Shi et al., 2019; Lu et al., 2022b) proposed to retrieve images corresponding to texts from the image corpus and use visual knowledge to improve the performance on the downstream tasks. Recently, some researchers (Long et al., 2021; Yang et al., 2021; Zhu et al., 2022) proposed to utilize the powerful text-to-image technical to

obtain the imagination representation of language and infuse them into the language model via the prefix-tuning (Li and Liang, 2021) way. In this paper, we also compared the visual prefix-based prompt learning methods (Liang et al., 2022; Jin et al., 2022; Tsimpoukelli et al., 2021), which has been verified to improve the performance of pretrained language models.

3 Methodology

3.1 Overview

ModICR focuses on infusing the given multi-modal information: premise, image, and answer, into the language model to make conditional inferences based on textual and visual clues. The overview of ModICR is illustrated in Figure 2. Specifically, given the premise $P = (p_1, \dots, p_M)$, image I and answer candidates $A = (a_1, \dots, a_Y)$, where p_i, a_i indicate the i th token of premise and the i th answer in the candidate set respectively, we first use the visual encoder to obtain the image representation, which is projected into the visual prefix to provide the objective environment information. Considering a semantic gap between visual prefixes and text when the language model performs context learning, we devise an alignment mapping network based on a multi-grained vision-language semantic aligner to gain the cross-modal align-

ment prefix. Finally, the two-type prefixes, premise text, and answer candidate are fed to the language model via the instruction learning way to perform multi-modal context reasoning.

3.2 Base Model

Previous methods (Dong et al., 2022; Chen et al., 2020; Yu et al., 2021a) adopt the pretrained vision-language model to obtain joint representation of text and image during inferring. Similarly, we utilize the pretrained single-stream bidirectional encoder Oscar (Li et al., 2020) as the backbone of the visual encoder and multi-grained vision-language semantic aligner. In this case, the image feature is first extracted by the widely-used tool FasterRCNN (Ren et al., 2015) and fed into the visual encoder and aligner. Oscar mainly make the token-level semantic alignment between image and text. Hence, following Yang et al. (2022), we pre-train Oscar-based chunk-aware semantic interactor on the Flickr30k Entities (Plummer et al., 2015) data set to perform the phrase-level semantic alignment between text and image.

3.3 Mapping Networks

We denote the obtained sequence representation of the image and the text aligned with the image features to $\mathbf{H}_I = (\mathbf{h}_{I_g}, \mathbf{h}_{I_1}, \dots, \mathbf{h}_{I_O})$, $\mathbf{H}_{ta} = (\mathbf{h}_{tag}, \mathbf{h}_{ta_1}, \dots, \mathbf{h}_{ta_N})$, and $\mathbf{H}_{pa} = (\mathbf{h}_{pag}, \mathbf{h}_{pa_1}, \dots, \mathbf{h}_{pa_N})$, respectively, where \mathbf{h}_{I_i} indicates the output hidden state of i th image region (obtained by FasterRCNN). \mathbf{h}_{ta_i} or \mathbf{h}_{pa_i} represents the token-level or phrase-level aligned representation of i th token in answer text. N is the token length of answer. Similarly, \mathbf{h}_{I_g} , \mathbf{h}_{tag} , and \mathbf{h}_{pag} show the global representations of image, token-level and phrase-level alignment information, respectively. However, the obtained visual and alignment embedding vectors may lie in a representation space different from the language model (used in the multi-modal context reasoner) due to the discrepancy across models. To alleviate this gap, we adopt the feature mapping network (Mokady et al., 2021) to project them into the corresponding learnable prefixes.

Vision Mapping Network (VMN). As the top blue part shown in Figure 2, we use the visual encoder to encode the image and employ a vision mapping network to project image representation \mathbf{H}_I into the sequence of visual prefix $\mathbf{V} = (v_1, \dots, v_l)$ with the mixed length l . v_i represents the i th visual embedding. The workflow is

$$v_1, \dots, v_l = \text{VMN}(\mathbf{h}_{I_g}). \quad (1)$$

For VMN, we adopt a two-layer perceptron with a ReLU activation function. It could be pretrained on large-scale image-text pairs for projecting visual features into the visual prefix that has the same space distribution as word embedding in LMs.

Alignment Mapping Network (AMN). It is capable of capturing the multi-view semantic alignment information of image-text pair and converting it into the cross-modal alignment prefix. Such prefix can bridge the semantic gap between visual prefix and text in the language model, enhancing the interactive understanding of image-text information. Specifically, we first apply a two-layer transformer to capture the pivotal multi-view alignment information lied in \mathbf{H}_{ta} and \mathbf{H}_{pa} . The specific calculation process of the first layer is as follows:

$$\begin{aligned} \mathbf{h}_{dr} &= \mathbf{W}^{dr}([\mathbf{h}_{tag}, \mathbf{h}_{pag}]) + \mathbf{b}^{dr}, \\ \mathbf{h}_{cr} &= \text{cross}(\mathbf{h}_{dr}, [\mathbf{h}_{ta_1}, \dots, \mathbf{h}_{ta_N}, \mathbf{h}_{pa_1}, \dots, \mathbf{h}_{pa_N}]), \\ \mathbf{h}_{ag}^1 &= \text{MLP}(\mathbf{h}_{cr}), \end{aligned} \quad (2)$$

where \mathbf{W}^{dr} and \mathbf{b}^{dr} are learnable parameters. *cross* represents the cross-attention calculation process. $[\]$ shows the concatenate computation. After doing the same two-layer calculation, we obtain the pivotal alignment representation \mathbf{h}_{ag} . Secondly, we project it into the cross-modal alignment prefix via a similar calculation process as the vision mapping network (Eq. 1). Finally, we gain an alignment prefix representation $\mathbf{A} = (a_1, \dots, a_m)$, where a_i indicates the i th alignment embedding and m is the length of prefix. By doing so, AMN could capture the pivotal semantic alignment information and project them into the learnable prefix vectors in the word embedding space.

3.4 Multi-Modal Context Reasoner

After obtaining two types of the prefix, we infuse them into an context reasoner to conduct cross modal reasoning, where we adopt the pretrained language model RoBERTa (Liu et al., 2019) as the context reasoner. We utilize the widely used instruction-learning method to incorporate the whole context encoding information. Specifically, we fill visual prefix, alignment prefix, premise and answer candidate in a pre-defined instruction template, “<cls> Is Answer correct or wrong based on conditions? <sep> Conditions: The Image is <V>, Bridge between the following text and image

is $\langle \mathbf{A} \rangle$, *Premise Text* is $\langle \text{Premise Text} \rangle \langle \text{sep} \rangle$, *Answer* is $\langle \text{Answer candidate} \rangle$. These special symbols, $\langle \mathbf{V} \rangle$, $\langle \mathbf{A} \rangle$, $\langle \text{Premise Text} \rangle$, and $\langle \text{Answer candidate} \rangle$, will be replaced by the obtained prefix vectors \mathbf{V} and \mathbf{A} , and word embedding representations of premise and answer in turn. The sequence representation is fed into the context reasoner to infer the final result. This way, we can utilize the context learning capability of pretrained language model to tackle the multi-modal reasoning problem. We obtain the inferring result of each answer candidate by applying a two-layer perceptron with the ReLU activation function on the output hidden state \mathbf{h}_{cls} of the top layer in RoBERTa. The whole training objective of ModICR can be defined as

$$\mathcal{L}_f = - \sum_{i=1}^4 \log P_i(x_i = q), \quad (3)$$

where x_i is the output probability on i th answer candidate and q is the label.

3.5 Training and Inference

To make Eq. 2 in the alignment mapping network capture pivotal multi-view alignment information, we will first train it about one epoch for alleviating the cold start problem leading to the collapse of the network. Concretely, we use a linear function to project \mathbf{h}_{ag} into the confidence score and employ the cross entropy loss to optimize it locally with the golden label q . The training process is regarded as \mathcal{L}_1 . Thus, the whole training process could be defined as

$$\mathcal{L} = \begin{cases} \mathcal{L}_1, & \text{steps} < N_{\text{whole}}, \\ \mathcal{L}_f, & \text{steps} > N_{\text{whole}}, \end{cases}$$

where *steps* shows the optimization step during training and N_{whole} represents the start of the whole training.

For inference, we input each answer candidate with premise and image into ModICR to obtain the confidence score and adopt the maximum one as the final result.

4 Experiment

4.1 Data sets

Conditional inference on joint textual and visual clues is a task that the text provides the prior permutation or the complementary information (external knowledge) with the image. There are few data sets

that meet the above requirement in the community. To verify the effectiveness of the proposed model, we first adopt the high-quality human-constructed PMR (Dong et al., 2022) data set, which contains 12,080 training samples, 1,538 validation samples and 1,742 testing samples. Textual premises pass the human cross-check annotation and contain six categories: relationship, personality, mood, and so on. In addition, we also reorganized a corresponding large-scale data set according to the VCR data set (Zellers et al., 2019a). We combine the given correct rationale and question as the textual premise and reform the original task into inferring the answer based on the new premise and image, i.e., QR→A. This way, the rationale could provide external knowledge information different from the source image. We set the original validation as the test set and selected some training samples as the validation set. Finally, the samples are divided into 210k training/2,923 validating/ 26,534 testing.

4.2 Baselines

We compare the proposed method to pretrained LMs and VLMs as follows:

BERT (Devlin et al., 2019) and **RoBERTa** (Liu et al., 2019) are both the transformer-based large language model, having achieved impressive performance on many natural language understanding tasks. We fine-tune them with only access to the textual premise.

VL-BERT (Lu et al., 2019) is a dual-stream pretrained cross-modal model. It adopts the BERT architecture, and the visual feature are concatenated with text embedding.

ERNIE-VL (Yu et al., 2021a) is a single-stream fusion encoder. It utilizes the structured knowledge obtained from scene graphs to learn joint representations of vision and language.

UNITER (Chen et al., 2020) also expands the BERT architecture to incorporate visual information and power heterogeneous downstream vision-language tasks with joint multi-modal embeddings.

Oscar (Li et al., 2020) is also a single-stream fusion encoder that uses object tags detected in images as anchor points to ease the learning of alignments significantly.

OFA (Wang et al., 2022) is a sequence-sequence cross-modal learning framework that unifies a diverse set of cross-modal and unimodal tasks, including visual grounding, image captioning, image classification, language modelling, etc.

Method ↓ Types →	Validation	Testing
BERT-B (Devlin et al., 2019)	-	65.2
VL-BERT-B (Lu et al., 2019)	-	75.4
ERNIE-VL-B (Yu et al., 2021a)	-	79.0
UNITER-B (Chen et al., 2020)	-	77.4
Oscar-B (Li et al., 2020)	77.7	76.1
RoBERTa-L (Liu et al., 2019)	77.3	75.0
PromptFuse (Liang et al., 2022)	77.4	76.5
VL-BERT-L (Lu et al., 2019)	-	79.3
ERNIE-VL-L (Yu et al., 2021a)	-	<u>79.9</u>
UNITER-L (Chen et al., 2020)	-	77.0
OFA-L (Wang et al., 2022)	79.9	79.1
MVPTR (Li et al., 2022b)	79.5	78.9
CALeC (Yang et al., 2022)	<u>80.1</u>	78.7
ModCR (frozen VLMs)	85.0	84.3
ModCR (fine-tune VLMs)	85.8	84.7

Table 1: Model performance (accuracy) on the PMR data set. The results of BERT, VL-BERT, ERNIE-VL, and UNITER are reported by Dong et al. (2022). For baselines, “-B” and “-L” indicate the base and large version, respectively. The underscore and bold indicate the second highest value and best performance (same as following tables). “frozen VLMs” and “fine-tune VLMs” represent whether the parameters of the visual encoder and multi-grained vision-language aligner are involved in training.

MVPTR (Li et al., 2022b) is a pretrained cross model that introduces the multi-level semantic alignment of vision-language to facilitate representation learning synergistically.

CALeC (Yang et al., 2022) is a unified prediction and generation model for some vision-language tasks, which introduces the chunk-aware semantic interactor to improve the semantic alignment representation and uses the lexical constraint technical to promote the quality of generation.

PromptFuse (Liang et al., 2022) is a prompt-based learning method to infuse visual information into the language model. It randomly initializes two learnable vectors as the alignment prefix to improve the space representation projection of image and text and bridge the semantic gap between the visual prefix and text.

4.3 Implementation Details

We use the Adam (Kingma and Ba, 2014) optimizer to train the above models on 2 A100 GPUs with a base learning rate of $2e-5$, a batch size of 32, and a dropout rate of 0.1. For each sample, we set the maximum number of visual regions extracted by FasterRCNN to 10. We set N_{whole} to 1 epoch and

Method ↓ Types →	AT ↑	D1 ↓	AF ↓	D2 ↓
BERT-B (Devlin et al., 2019)	65.2	19.8	19.6	4.5
Oscar-B (Li et al., 2020)	76.1	10.2	12.1	1.7
RoBERTa-L (Liu et al., 2019)	75.0	17.7	6.1	1.2
PromptFuse (Liang et al., 2022)	76.5	16.5	5.8	1.2
ERNIE-VL-L (Yu et al., 2021a)	79.9	10.7	8.2	1.2
OFA-L (Wang et al., 2022)	79.1	9.7	9.9	1.3
MVPTR (Li et al., 2022b)	78.9	7.5	11.8	1.8
CALeC (Yang et al., 2022)	78.7	8.6	10.9	1.8
ModCR (frozen VLMs)	84.3	9.2	5.6	0.9
ModCR (fine-tune VLMs)	84.7	7.8	6.8	0.7

Table 2: Detailed performance of models on the test set of PMR. The results of BERT and ERNIE-VL are reported by Dong et al. (2022). AT, D1, AF, D2 represent the Action True and Image True, Action True yet Image False, Action False yet Image True, Action False and Image False, respectively. “Action True or False” indicate the answer whether meets the premise. Similarly, “Image True or False” show the answer whether meets the image information.

Method ↓ Types →	Validation	Testing
Oscar-B (Li et al., 2020)	87.3	86.0
RoBERTa-L (Liu et al., 2019)	<u>92.7</u>	<u>91.8</u>
OFA-L (Wang et al., 2022)	90.3	89.4
MVPTR (Li et al., 2022b)	84.2	85.3
CALeC (Yang et al., 2022)	90.8	90.5
ModCR (frozen VLMs)	94.5	93.6
ModCR (fine-tune VLMs)	94.7	94.0

Table 3: Model performance (accuracy) on the validation and testing sets of VCR (QR→A) data set.

adopt the pre-trained parameters of the base version of Oscar to initialize the multi-grained vision-language semantic aligner. While training the chunk-level semantic interactor on the Flickr30k Entities data set, we follow the parameter settings presented in Yang et al. (2022) and train it for about ten epochs. We adopt the Roberta_{large} to initialize the multi-modal context reasoner. The visual and cross-modal alignment prefix lengths are both set to 5. All methods performed on the two data sets employ the validation set to select the best-performing model.

4.4 Main Results

Overall Performance. We report the performance of models on PMR and VCR (QR→A) data sets, which are shown in Tables 1 and 3. From the whole experimental results, we observe that the proposed method significantly outperforms previously strong

Method ↓ Types →	Validation	Testing
CALeC (Yang et al., 2022)	80.1	78.7
RoBERTa-L (Liu et al., 2019)	77.3	75.0
PromptFuse (LV=1, LA=2)	77.4	76.5
ModCR (LV=1, LA=0)	78.1	76.0
ModCR (LV=3, LA=0)	78.2	77.8
ModCR (LV=5, LA=0)	77.3	76.8
ModCR (LV=3, LA=1)	84.9	83.5
ModCR (LV=3, LA=5)	85.8	83.9
ModCR (LV=3, LA=7)	85.3	84.1
ModCR (LV=1, LA=1)	84.0	82.3
ModCR (LV=3, LA=3)	84.8	83.8
ModCR (LV=5, LA=5)	85.0	84.3
ModCR (LV=7, LA=7)	85.1	82.8
ModCR (LV=10, LA=10)	79.7	79.3

Table 4: The experimental results of ModCR with different prefix length on the PMR data set. We frozen the parameters of VLMs for all ModCR variants. “LV” and “LA” indicate the lengths of visual and alignment prefix respectively, where “=0” represents that the corresponding mapping network is removed.

baselines such as gain by 5.7%, 4.8% on the validation and testing of the PMR data set compared to CALeC and ERNIE-VL-L. According to the performance of BERT-B and RoBERTa (only text input), we know that the premise can provide vital information to infer the correct option. The performance is further improved when combined with visual content and cross-modal semantic alignment prefix for inference, e.g., ModCR (frozen VLMs) vs. RoBERTa: 84.3 vs. 75.0, PromptFuse vs. RoBERTa: 76.5 vs. 75.0. For model performances on VCR (QR-A), however, we observe that the pretrained VLMs have worse performance compared to RoBERTa-L, which displays that VLMs do not make good use of the abstract semantics of the premise for contextual reasoning. ModCR that takes the RoBERTa-L as the main backbone surpasses pretrained VLMs and LMs on two data sets, which suggests that our method effectively utilizes the semantic information of different modalities while performing reasoning.

Is Context Reasoning Capability Improved?

We present the detailed performances of models on the test set of PMR to check the ability of models to infer different types of answer candidates, which contain AT, D1, AF, and D2, as shown in Table 2. The reported results indicate that RoBERTa better uses the abstract semantic information of

MappNet	RoBERTa	VLM	Validation	Testing
✓	×	×	85.7	85.8
✓	✓	×	94.5	93.6
✓	✓	✓	94.7	94.0
✓	×	×	72.2	69.2
✓	✓	×	85.0	84.3
✓	✓	✓	85.8	84.7

Table 5: The detailed performance of ModCR with different training strategies. “MappNet” indicates the two types of mapping networks. “✓” represents parameters of the module will be updated during training. The top 3 lines show the experimental results on the VCR (QR→A), and the bottom 3 lines is PMR.

premise to infer the correctness of the following action compared to VLMs, e.g., RoBERTa without visual information has the lowest error rate across all baselines in action recognition (AT). In addition, we also find that although the ability of recently proposed VLMs to reason with abstract textual clues has been improved, there is still a particular gap compared to LMs, e.g., AT performance: OFA-L (8.2) vs. RoBERTa (6.0). When employing the language model RoBERTa as the reasoner and infusing the visual information in it, we observe that the overall accuracy of the model is further improved. However, the previous vision-infusing method has a low utilization rate of visual information (D1: 16.5 for PromptFuse). As the bottom two lines shown in Table 2, ModCR, which utilizes the multi-view text-image semantic alignment information, maintains the abstract reasoning ability based on premise and also substantially improves the utilization rate of image information.

Through the above analysis, we can obtain that it is necessary to introduce vision-language semantic alignment information for vision-aided language models. Furthermore, there is still a large room for improvement in the contextual reasoning capability of the pretrained VLMs.

4.5 Ablation Studies

To analyze the effectiveness of ModCR in detail, we design multiple model variants and the experimental results are shown in Tables 4 and 5. We select the high-quality PMR (manual annotation and inspection) data set as the experimental scene of ablation studies. For PromptFuse (Liang et al., 2022), we adopt RoBERTa-L as the backbone and all parameters are updated during training.

Is Alignment Mapping Network Effective?

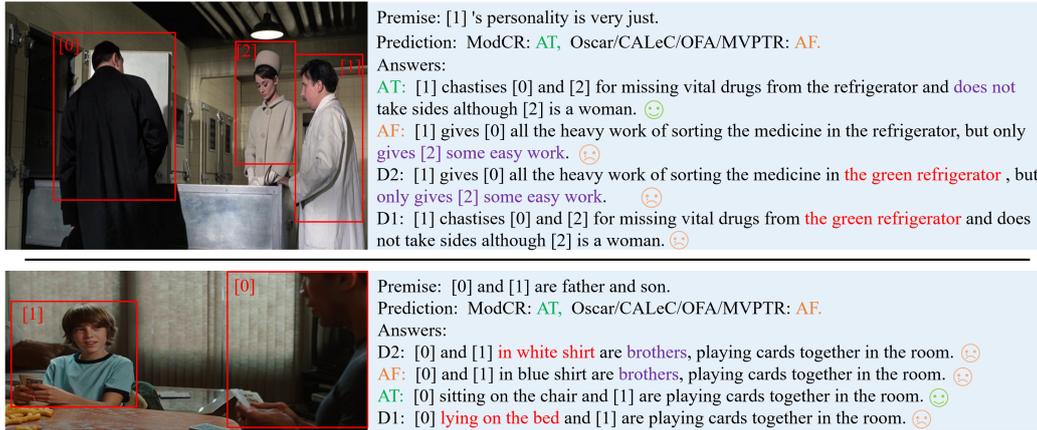


Figure 3: Two cases from the test set of PMR. Different persons are represented by red squares and numbers in the left images. For answer candidates, the red words indicate the content that does not meet the image; the purple indicates the content that is conflict with the premise text clues; and the green and orange indicate the correct option and the wrong answer respectively.

From Table 4, comparing ModCR performances with $LA=0$ and $LA \geq 1$, we observe that the performance of ModCR drops markedly when it abandons vision-language semantic alignment information. Compared to PromptFuse that randomly initializes two learnable alignment prefix vectors, the proposed alignment mapping network equipped with the multi-grained cross-modal aligner is more effective, e.g., PromptFuse vs. RoBERTa-L: 76.5 vs. 75.0, and performance comparisons of ModCR vs. RoBERTa-L.

Effect of Prefix Length on Model Performance.

From the performance of the visual prefix and alignment prefix at different lengths in Table 4, we can see that the performance of ModCR varies greatly under different lengths for the two types of prefix. The ModCR performs best when both prefixes are taken as 5. Furthermore, excessively long visual prefixes impair the overall performance, which may be attributed to the fact that redundant and inaccurate visual prefix has an inferior effect on the context learning capability of language model.

Model Performance with Different Training Strategies.

We present the detailed performance of ModCR with different training strategies on Table 5. By comparing the experimental results of “frozen VLM” and “fine-tune VLM” on two data sets, we observe that the performance of the proposed method is further improved when all parameters of ModCR are updated during training. Although the training speed is slower, this could further integrate the complementary reasoning capabilities of VLM and LM. In addition, only finetuning

MappNet has inferior performances, which may be addressed via pretraining on external large-scale image-text corpus.

4.6 Case Study

We report two cases in Figure 3 to analyse the performance of models in detail. The premise texts of two samples are about the character (top case) and relationship (bottom one) of persons respectively. Although pre-trained VLMs can infer whether the answer candidate satisfies the image content, they cannot effectively use the premise information to perform reasoning. Contrastly, ModCR utilizes the two-modal semantic information to determine the correct answer. It indicates that regrading two different cues as pre-context states and employing the context reasoning ability of language models is a simple and effective approach for cross modal reasoning tasks. In addition, ModCR could infer the description “in white shirt” and “lying on the bed” do not meet the image content (the boy wearing blue shirt and sitting on the chair), which may be attributed to the semantic aligner. To conclude, the alignment prefix can improve the whole performance of allowing the language model to understand the visual information and perform reasoning.

5 Conclusion and Future Work

In this paper, we propose a multi-modal context reasoning approach named ModCR for the scenario of conditional inference on joint visual and textual clues. It regards the given image and text as the two

types of pre-context states and infuses them into the language model via the instruction learning method to perform such multi-modal reasoning. The experimental results on two data sets show the effectiveness of ModCR. For the future, we will explore two research directions: 1) how to improve the context learning capability of pretrained VLMs. 2) exploring the conditional inference on complex visual and textual clues, where it contains multiple clues lying in more modalities.

Limitations

The proposed method has several limitations: 1) The current approach achieves hunky context reasoning performance in the cross-modal scene of a single text clue and image, but the context reasoning capability in the scene containing multiple textual and visual clues still needs to be further explored, such as video and long text. 2) From the experimental results, we observed that the visual prefix length greatly impacts the stability of language models infused with visual information. Hence, we still need to explore effective and stable vision-aided language models for natural language processing and multi-modal scenarios. 3) We also hope this work could spark further research on improving the long context reasoning capability of pretrained vision-language models.

References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Regina Barzilay and Mirella Lapata. 2005. Collective content selection for concept-to-text generation. In *Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, pages 331–338.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. 2021. [Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language BERTs](#). *Transactions of the Association for Computational Linguistics*, 9:978–994.
- Cheng Chen, Zhenshan Tan, Qingrong Cheng, Xin Jiang, Qun Liu, Yudong Zhu, and Xiaodong Gu. 2022. Utc: A unified transformer with inter-task contrastive learning for visual dialog. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18103–18112.
- Hongyu Chen, Ruifang Liu, and Bo Peng. 2021. Cross-modal relational reasoning network for visual question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3956–3965.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *ECCV*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2020. e-snli-ve: Corrected visual-textual entailment with natural language explanations. *arXiv preprint arXiv:2004.03744*.
- Qingxiu Dong, Ziwei Qin, Heming Xia, Tian Feng, Shoujie Tong, Haoran Meng, Lin Xu, Zhongyu Wei, Weidong Zhan, Baobao Chang, et al. 2022. Premise-based multimodal reasoning: Conditional inference on joint textual and visual clues. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 932–946.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Feng Gao, Qing Ping, Govind Thattai, Aishwarya Reganti, Ying Nian Wu, and Prem Nataraajan. 2022. Transform-retrieve-generate: Natural language-centric outside-knowledge visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5067–5077.
- Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239.

- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR.
- Woojeong Jin, Yu Cheng, Yelong Shen, Weizhu Chen, and Xiang Ren. 2022. [A good prompt is worth millions of parameters: Low-resource prompt-based learning for vision-language models](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2763–2775, Dublin, Ireland. Association for Computational Linguistics.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Benno Kroger, Vaibhav Adlakha, Vibhav Vineet, Yash Goyal, Edoardo Ponti, and Siva Reddy. 2022. Image retrieval from contextual descriptions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3426–3440, Dublin, Ireland. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Jiangtong Li, Li Niu, and Liqing Zhang. 2022a. From representation to reasoning: Towards both evidence and commonsense reasoning for video question-answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21273–21282.
- Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. 2021. Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. *ACL*.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV 2020*.
- Zejun Li, Zhihao Fan, Huaixiao Tou, Jingjing Chen, Zhongyu Wei, and Xuanjing Huang. 2022b. Mvptr: Multi-level semantic alignment for vision-language pre-training via multi-stage learning. In *Proceedings of the 30th ACM International Conference on Multimedia, MM '22*, page 4395–4405, New York, NY, USA. Association for Computing Machinery.
- Sheng Liang, Mengjie Zhao, and Hinrich Schuetze. 2022. [Modular and parameter-efficient multimodal fusion with prompting](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2976–2985, Dublin, Ireland. Association for Computational Linguistics.
- Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. 2021. M6: A chinese multimodal pretrainer. *arXiv preprint arXiv:2103.00823*.
- Weizhe Lin and Bill Byrne. 2022. [Retrieval augmented visual question answering with outside knowledge](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11238–11254, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 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*.
- Quanyu Long, Mingxuan Wang, and Lei Li. 2021. [Generative imagination elevates machine translation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5738–5748, Online. Association for Computational Linguistics.
- Siqu Long, Feiqi Cao, Soyeon Caren Han, and Haiqing Yang. 2022. Vision-and-language pretrained models: A survey. *IJCAI*.
- Haoyu Lu, Nanyi Fei, Yuqi Huo, Yizhao Gao, Zhiwu Lu, and Ji-Rong Wen. 2022a. Cots: Collaborative two-stream vision-language pre-training model for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15692–15701.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Yujie Lu, Wanrong Zhu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2022b. Imagination-augmented natural language understanding. *NACCL*.
- Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*.

- Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. [Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models](#). In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2641–2649.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. [Faster r-cnn: Towards real-time object detection with region proposal networks](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. 2019. Kvqa: Knowledge-aware visual question answering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8876–8884.
- Haoyue Shi, Jiayuan Mao, Kevin Gimpel, and Karen Livescu. 2019. [Visually grounded neural syntax acquisition](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1842–1861, Florence, Italy. Association for Computational Linguistics.
- Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. [OFA: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework](#). In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 23318–23340. PMLR.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*.
- Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony Dick, and Anton Van Den Hengel. 2017. Visual question answering: A survey of methods and datasets. *Computer Vision and Image Understanding*, 163:21–40.
- Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. 2019. Visual entailment: A novel task for fine-grained image understanding. *arXiv preprint arXiv:1901.06706*.
- Qian Yang, Yunxin Li, Baotian Hu, Lin Ma, Yuxin Ding, and Min Zhang. 2022. [Chunk-aware alignment and lexical constraint for visual entailment with natural language explanations](#). In *Proceedings of the 30th ACM International Conference on Multimedia*, MM '22, page 3587–3597, New York, NY, USA. Association for Computing Machinery.
- Ze Yang, Wei Wu, Huang Hu, Can Xu, Wei Wang, and Zhoujun Li. 2021. Open domain dialogue generation with latent images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14239–14247.
- Keren Ye and Adriana Kovashka. 2021. A case study of the shortcut effects in visual commonsense reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 3181–3189.
- Fei Yu, Jiji Tang, Weichong Yin, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2021a. Ernie-vil: Knowledge enhanced vision-language representations through scene graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3208–3216.
- Youngjae Yu, Jiwan Chung, Heeseung Yun, Jongseok Kim, and Gunhee Kim. 2021b. Transitional adaptation of pretrained models for visual storytelling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12658–12668.
- Abdulganiyu Abdu Yusuf, Feng Chong, and Mao Xi-anling. 2022. An analysis of graph convolutional networks and recent datasets for visual question answering. *Artificial Intelligence Review*, pages 1–24.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019a. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6720–6731.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019b. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Shunyu Zhang, Xiaoze Jiang, Zequn Yang, Tao Wan, and Zengchang Qin. 2022. Reasoning with multi-structure commonsense knowledge in visual dialogue. In *Proceedings of the IEEE/CVF Conference*

on Computer Vision and Pattern Recognition, pages 4600–4609.

Wanrong Zhu, An Yan, Yujie Lu, Wenda Xu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2022. Visualize before you write: Imagination-guided open-ended text generation. *arXiv preprint arXiv:2210.03765*.