# Coarse-to-Fine Reasoning for Visual Question Answering

Binh X. Nguyen[1], Tuong Do[1], Huy Tran[1], Erman Tjiputra[1], Quang D. Tran[1], Anh Nguyen[2]

[1]AIOZ, Singapore

[2]University of Liverpool, UK

{binh.xuan.nguyen, tuong.khanh-long.do, huy.tran, erman.tjiputra,quang.tran}@aioz.io

anh.nguyen@liverpool.ac.uk

## Abstract

*Bridging the semantic gap between image and question is an important step to improve the accuracy of the Visual Question Answering (VQA) task. However, most of the existing VQA methods focus on attention mechanisms or visual relations for reasoning the answer, while the features at different semantic levels are not fully utilized. In this paper, we present a new reasoning framework to fill the gap between visual features and semantic clues in the VQA task. Our method first extracts the features and predicates from the image and question. We then propose a new reasoning framework to effectively jointly learn these features and predicates in a coarse-to-fine manner. The intensively experimental results on three large-scale VQA datasets show that our proposed approach achieves superior accuracy comparing with other state-of-the-art methods. Furthermore, our reasoning framework also provides an explainable way to understand the decision of the deep neural network when predicting the answer. Our source codes can be found at: https://github.com/aioz-ai/CFR_VQA*

## 1. Introduction

The Visual Question Answering (VQA) task aims to predict the correct answer of a given question such that the answer is consistent with the visual image content. There are two main variants of VQA, i.e., Free-Form Opened-Ended (FFOE) and Multiple Choice (MC). In FFOE, an answer is a free-form response of a given image-question input pair, while in MC, the answer is chosen from a list of predefined ground-truth. In both cases, extracting meaningful features from the images and questions plays a key role. Furthermore, mapping the semantic features from the images and questions also strongly affects the results [16]. Most of the existing solutions for the VQA task rely on visual relations [4, 5, 56, 60], attention mechanisms [26, 47, 50], external knowledge [18, 29], or message passing [50] to link the visual clue with the associated information in the question.

While both extracting and reasoning the features of the image and question are important for VQA, they are not trivial tasks in practice. Many questions (and answers) are composed of complex semantic information, which can have noise or ambiguous attributes. Current methods focus on utilizing visual information [5,22,33,35,38,40,42,44,53] without considering if the supporting information is useful or not [16]. Besides, many approaches aim to enrich the information extracted from both image and question regardless of the noisy information that may occur [5, 14, 15, 18]. This leads to the fact that although the image and question features can be extracted by a deep convolutional neural network, they may not be effectively utilized to reason and predict the correct answer.

To bridge the semantic gap between images and questions in VQA, we introduce a new framework that focuses on reasoning the visual contents in the image and the semantic clues in the question in a coarse-to-fine manner. Our observation is that both image and question's features can be extracted gradually at different fine-grained levels. Therefore, we can map these features in each level to allow a stronger connection when reasoning. Our framework contains effective extractors for extracting meaningful features and predicates from the image and question. Furthermore, the answer outputted by our framework can be reasoned explicitly through the distribution maps during the prediction progress. These maps indicate the necessity of input features or predicates, allow us to understand which information is meaningful for predicting the answer. Our contributions can be summarized as follows:

- We propose a simple, yet effective framework to extract meaningful features and predicates from the question and image. The extracted information can be used to explain the decision of the deep network.

- We introduce a new coarse-to-fine reasoning method to bridge the semantic gap between the question and

image when predicting the answer.

- We conduct intensive experiments to validate our method. Our source code and trained models will be released for further study.

## 2. Related Work

There are numerous reasoning VQA methods [1, 6, 12, 14, 19, 36, 37, 39, 48, 52, 55, 58, 63, 64] that focus on learning the relations between visual regions and words in questions implicitly, e.g., through message passing [50], pairwise relationship modeling [4], adversarial learning [8, 31, 51], or graph parsing methods defined by inter/intra-class edges [15]. Other works focus on leveraging external information [18] or explicit scene graph [5] to extract features from input images. ReGAT [30] considers both explicit and implicit relations to enrich image representations. Most of the current VQA works focus on enriching image representation without examining whether the enriched information is necessary for reasoning the answer or not [28].

Extracting meaningful features from images, questions, and their joint embedding is crucial in the VQA task. For image representation, grid features [23, 65] or object features [2, 13, 43, 46, 49] are widely used. For question embedding, Glove [45] and BERT [10] are used to present words and sentences. Besides, using large-scale pre-training models on image-text pairs is also popular [7, 32]. For learning the joint embedding, many approaches use attention mechanisms [11, 26, 41, 47, 50, 61, 62]. The authors in [57] propose Stacked Attention Networks to localize image regions that are relevant to the question. In [26], the authors propose Bilinear Attention Networks for VQA. Recently, in [11], the authors introduce Compact Trilinear Interaction which simultaneously learns the interaction between images, questions, and answers.

Unlike other approaches that focus on enriching information from image and question, in this work, we consider the interaction among the semantic clues in questions and the visual contents of the image ranging from object-level to fine-grained level. Hence, we apply a simplified fine-grained detector inspired by Faster R-CNN model [46] to extract visual features and predicates, rather than leveraging complicated scene graph generators. This setup allows us to achieve competitive results compared with other approaches, while keeping the network at a reasonable computational cost.

## 3. Methodology

### 3.1. Overview

Our Coarse-to-Fine Reasoning (CFR) framework takes an image and a question as inputs. The image is passed through the Image Embedding module to extract the re-

gion of interest (RoI) features and visual predicates. The question is processed in the Question Embedding module to extract the question features and question predicates. The predicates are keywords about objects, relations, or attributes of the image/question. To effectively map the visual modality and language modality, we jointly learn their features, as well as their predicates in the Coarse-to-Fine Reasoning module. Figure 1 illustrates an overview of our framework.

### 3.2. Image Embedding

The goal of the Image Embedding module is to extract RoI features and visual predicates from the input image. The RoI features are extracted by a deep object detector to localize all potential regions of interest. The visual predicates are extracted by classifying attributes and relations based on the visual RoI features provided by the object detector.

In practice, as in [26, 49], we use the pre-trained Faster R-CNN model [2] to extract visual features for each RoI. Note that the RoI feature is an important visual input for the VQA task. Therefore, we retain the original Faster R-CNN multi-task loss for object detection, then adding two additional Cross-Entropy losses for attribute class predictor and relation class predictor. The extracted objects, as well as their attributes and relations, are then re-arranged to form predicates. Each predicate follows one of three forms: single predicate <obj>; attribute-based predicate <attr, obj>; and relation-based predicate <obj1, rel, obj2>. Following [26, 30, 49], we use a pretrained Faster R-CNN model on the Visual Genome dataset [27] to extract predicates from the images. For each word in each predicate, we apply 300-dim Glove word embedding [45] to extract predicate features.

### 3.3. Question Embedding

The Question Embedding module aims to extract question features and question predicates. To extract question features, following [11, 26, 59], we apply 600-dim Glove word embedding [45] accompanied by GRU [9] to extract the features and learn the dependencies of all words in the question.

To extract question predicates, we pass the whole question through a stop-word filter. The filter is the combination of two lists. The first list contains words in the NLTK based stop-words [34] list, i.e., words that do not add much meaning in a sentence. The second list contains words from all the questions that have the frequency of occurrence is less than 10. Words in the second list are considered as rare words and hard for the model to learn. For each word in each question predicate, we apply 300-dim Glove word embedding [45] to extract the predicate features.
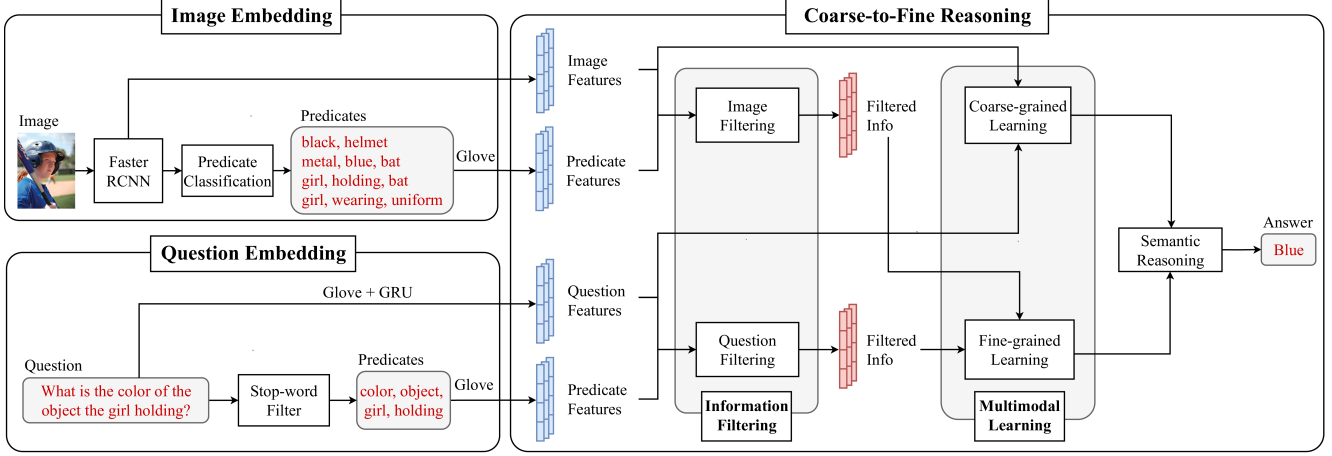
Figure 1. An overview of our framework.

## 3.4. Coarse-to-Fine Reasoning

Given the image features and predicates $(\mathbf{f}_i, \mathbf{p}_i)$ as well as the question features and predicates $(\mathbf{f}_q, \mathbf{p}_q)$, our goal is to predict an answer $\alpha$ in a list of ground-truth $\mathcal{A}$ using a trainable model $\theta$ as follow:

$$\hat{\alpha} = \underset{\alpha \in \mathcal{A}}{\arg\max}\, \theta\left(\alpha | \mathbf{f}_i, \mathbf{p}_i, \mathbf{f}_q, \mathbf{p}_q\right) \tag{1}$$

To effectively map the information of the question to the visual information in the image, the Coarse-to-Fine Reasoning module utilizes three steps: Information Filtering, Multimodal Learning, and Semantic Reasoning. The Information Filtering aims to filter out unnecessary visual information from the image based on the predicates. The Multimodal Learning module learns the semantic mapping between the question and image at coarse-grained and fine-grained levels. Finally, the Semantic Reasoning module combines the output of the multimodal learning step to predict the answer.

### 3.4.1 Information Filtering

Since the features and predicates of both the question and image are extracted by pretrained models, they may have noise or incorrect information. Therefore, we design the Information Filtering module to filter out unnecessary information. In practice, this module also helps us understand the importance of each RoI for each question. The Information Filtering takes the feature $\mathbf{f} \in \mathbb{R}^{n_f \times d_f}$ and the predicate $\mathbf{p} \in \mathbb{R}^{n_p \times d_p}$ as input. Both $\mathbf{f}$ and $\mathbf{p}$ have a matrix form; $n_f$, $n_p$ denote the number of instances (e.g., number of RoIs or number of predicates); $d_f$, $d_p$ denote the dimension of each instance.

To filter out the unnecessary information in the feature $\mathbf{f}$, we consider the predicate $\mathbf{p}$ as the supervision information.

Through the interaction mechanism, we compute a weighting map $\hat{\Psi} \in \mathbb{R}^{n_f}$ which is then applied to output the filtered information $\Psi \in \mathbb{R}^{n_f \times d_\Psi}$. $\hat{\Psi}$ is computed as follow:

$$\hat{\Psi} = \mathrm{softmax}\left(\sum_{i=1}^{n_p} \tau_f\left(\mathbf{f}\right) \tau_p\left(\mathbf{p}\right)_i^T\right) \tag{2}$$

where $\tau_f(\cdot)$ and $\tau_p(\cdot)$ are learnable linear projection funtions which project $\mathbf{f} \in \mathbb{R}^{n_f \times d_f}$ and $\mathbf{p} \in \mathbb{R}^{n_p \times d_p}$ into $\mathbf{f}' \in \mathbb{R}^{n_f \times d_\Psi}$ and $\mathbf{p}' \in \mathbb{R}^{n_p \times d_\Psi}$, respectively.

Given the weighting map $\hat{\Psi}$, the filtered information $\Psi$ is calculated by Equation (3):

$$\Psi = \left(\hat{\Psi} \cdot \mathbb{1}^T\right) \odot \tau_f\left(\mathbf{f}\right) + \tau_f\left(\mathbf{f}\right) \tag{3}$$

where $\mathbb{1} \in \mathbb{R}^{d_\Psi}$ is a channel-scaled vector; $\odot$ denotes the Hardamard product.

In practice, the Information Filtering module is applied on both the image features and predicates $(\mathbf{f}_i, \mathbf{p}_i)$, as well as the question features and predicates $(\mathbf{f}_q, \mathbf{p}_q)$ to achieve the filtered information $\Psi_i$ and $\Psi_q$. Here we use the unified symbol $\Psi$ for simplicity.

### 3.4.2 Multimodal Learning

Inspired by the Unitary Attention Mechanism [26], we design the Multimodal Learning module to jointly learn the features from the visual and language modalities. Multimodal learning is essential for identifying the correlation between each instance in the image and the question, then identifying which instances in the image are useful for answering the question.

In this module, the features are jointly learned at two levels: coarse-grained and fine-grained. The coarse-grained level learns the interaction between question features and
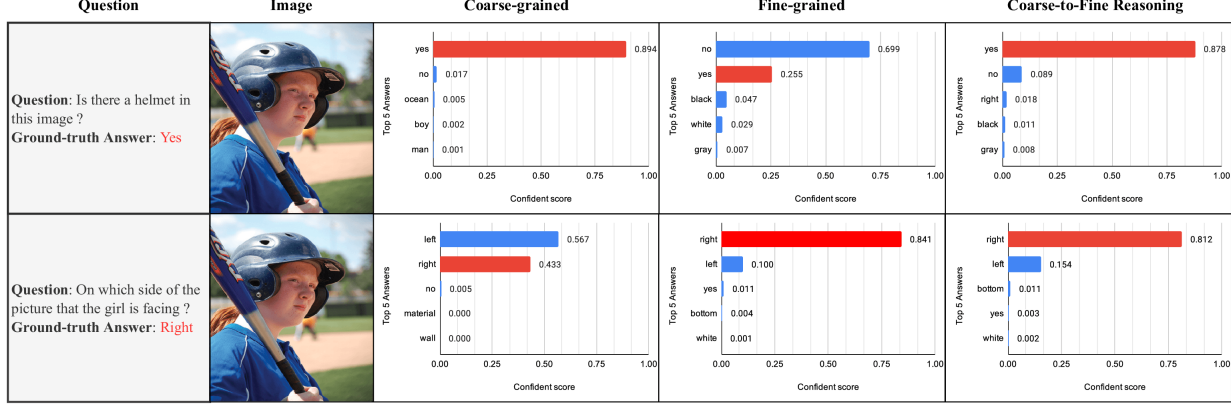
Figure 2. Examples of the predicted confidence scores of the Coarse-grained Learning, Fine-grained Learning, and Coarse-to-Fine Reasoning module.

image features, while the fine-grained level learns the interaction between filtered information of the image and question obtained from the Information Filtering step.

**Coarse-grained learning.** The inputs for coarse-grained learning are the image features $\mathbf{f}_i$ and question features $\mathbf{f}_q$. The output of coarse-grained learning is a joint representation $\mathbf{j}^{cg} \in \mathbb{R}^{d_{cg}}$, where $d_{cg}$ is the dimension of the joint representation. Each $k$-th element of the join representation $\mathbf{j}^{cg}$ is computed as follows:

$$\mathbf{j}_k^{cg} = (\mathbf{f}_q \mathbf{M}_{\mathbf{f}_q})_k^T \mathbf{A}^{cg} (\mathbf{f}_i \mathbf{M}_{\mathbf{f}_i})_k \tag{4}$$

where $\mathbf{M}_{\mathbf{f}_q} \in \mathbb{R}^{d_q \times d_{cg}}$ and $\mathbf{M}_{\mathbf{f}_i} \in \mathbb{R}^{d_i \times d_{cg}}$ are learnable factor matrices; $n_q$, $n_i$ denote the number of instances in question and image; $\mathbf{A}^{cg} \in \mathbb{R}^{n_q \times n_i}$ is the bilinear attention distribution map of the joint representation $\mathbf{j}^{cg}$; $d_q$, $d_i$ denote the dimension of each instance. The subscript $k$ in-

dicates the index of matrix column. $\mathbf{A}^{cg}$ is computed by Equation (5):

$$\mathbf{A}^{cg} = \mathrm{softmax}\left( \left( \mathbf{f}_q \mathbf{M}'_{\mathbf{f}_q} \right) \left( \mathbf{f}_i \mathbf{M}'_{\mathbf{f}_i} \right)^T \right) \tag{5}$$

where $\mathbf{M}'_{\mathbf{f}_q} \in \mathbb{R}^{d_q \times d_{cg}}$ and $\mathbf{M}'_{\mathbf{f}_i} \in \mathbb{R}^{d_i \times d_{cg}}$ are learnable factor matrices, and independent of $\mathbf{M}_{\mathbf{f}_q}$ and $\mathbf{M}_{\mathbf{f}_i}$.

**Fine-grained learning.** We apply the same process of coarse-grained learning for fine-grained learning. The only difference is the inputs for fine-grained learning are the image filtered information $\Psi_i$ and question filtered information $\Psi_q$. Similar to Equation 4 and 5, the fine-grained joint representation is computed as follow:

$$\mathbf{j}_k^{fg} = (\Psi_q \mathbf{M}_{\Psi_q})_k^T \mathbf{A}^{fg} (\Psi_i \mathbf{M}_{\Psi_i})_k \tag{6}$$

where $\mathbf{A}^{\text{fg}}$ is computed as:

$$\mathbf{A}^{\text{fg}} = \text{softmax}\left( \left(\Psi_{\text{q}}\mathbf{M}'_{\Psi_{\text{q}}}\right) \left(\Psi_{\text{i}}\mathbf{M}'_{\Psi_{\text{i}}}\right)^T \right) \qquad (7)$$

### 3.4.3 Semantic Reasoning

The goal of Semantic Reasoning is to selectively learn information from both the Coarse-grained and the Fine-grained learning steps using a learnable adaptive weight $\mathbf{W} \in \mathbb{R}^{|\mathcal{A}|}$, where $|\mathcal{A}|$ is the number of possible answers. In practice, this module takes $\mathbf{j}^{\text{cg}}$ and $\mathbf{j}^{\text{fg}}$ as inputs and then outputs the distribution $\rho \in \mathbb{R}^{|\mathcal{A}|}$ over candidates of all answers $\mathcal{A}$.

$$\rho = \text{softmax}\left(\mathbf{W}\tau\left(\mathbf{j}^{\text{cg}}\right) + \mathbf{W}'\tau'(\mathbf{j}^{\text{fg}})\right)$$
$$\text{s.t} \sum ||[\mathbf{W}_\alpha, \mathbf{W}'_\alpha]|| = 1, \forall \alpha \in \mathcal{A} \qquad (8)$$

where $\mathbf{W}$ and $\mathbf{W}'$ are the learnable adaptive weights of coarse-grained learning and fine-grained learning; $\tau(\cdot)$ and $\tau'(\cdot)$ are learnable projection functions that project $\mathbf{j}^{\text{cg}} \in \mathbb{R}^{\text{d}_{\text{cg}}}$ and $\mathbf{j}^{\text{fg}} \in \mathbb{R}^{\text{d}_{\text{fg}}}$ into $\rho^{\text{cg}} \in \mathbb{R}^{|\mathcal{A}|}$ and $\rho^{\text{fg}} \in \mathbb{R}^{|\mathcal{A}|}$, respectively. To satisfy the constraint in Equation (8), we apply the softmax function for each vector $[\mathbf{W}_\alpha, \mathbf{W}'_\alpha]$; the subscript $\alpha$ indicates the index of an answer in the answer list $\mathcal{A}$.

Through an end-to-end training process, the learned adaptive weights $\mathbf{W}$ identify the contribution of each input information to predict the answer. These weights are expected to robust with noisy information from the question or image at both the coarse-grained and the fine-grained level.

## 4. Experiments

### 4.1. Dataset, baseline and evaluation protocol

**Dataset.** We use three popular datasets in our experiments: GQA [21], VQA 2.0 [17], and Visual7W [65] . We follow the same split in each dataset for training and testing.

**Implementation.** We conduct experiments on an NVIDIA TITAN V 12GB GPU. The network is trained with a batch size of 32 and a learning rate of 0.001 using Adam optimizer. Following [25, 26, 49, 57], we use the Visual Genome [27] and Glove [45] to extract the image embedding and question embedding. Then we train the whole framework from scratch. The parameters $\text{d}_{\text{cg}}$ and $\text{d}_{\text{fg}}$ are empirically set to 768. The learnable factor matrices $\mathbf{M}$, $\mathbf{W}$ are initialized randomly at the beginning of the training phase and being learned through the training process. It takes approximately 10, 20, and 35 hours to train our network on Visual7W, VQA2.0, and GQA dataset, respectively.

**Baselines.** We compare our results with various recent methods in VQA. These methods can be categorized into three groups: joint learning mechanisms: BAN [26], Pythia [24], DFAF [15], fPMC [20], STL [54], CTI [11],

| Method | GQA (Acc) val | GQA (Acc) tes-dev | VQA 2.0 (Acc) val | VQA 2.0 (Acc) test-dev | Visual7W (Acc-MC) val | Visual7W (Acc-MC) test |
|---|---|---|---|---|---|---|
| BAN [26] | 61.5 | 55.2 | 66.0 | 70.0 | 65.7 | 67.5 |
| Pythia [24] | – | – | 66.3 | 70.0 | – | – |
| DFAF [15] | – | – | 66.2 | 70.2 | – | – |
| fPMC [20] | – | – | 61.7 | 63.9 | – | 66.0 |
| STL [54] | – | – | – | – | 67.5 | 68.2 |
| CTI [11] | 61.7 | 54.9 | 66.0 | 70.1 | 67.0 | 69.3 |
| MCAN [59] | – | 57.4 | 67.2 | 70.6 | – | – |
| MuRel [4] | – | – | 65.1 | 68.0 | – | – |
| ReGAT [30] | – | – | 67.2 | 70.3 | – | – |
| MMN [6] | – | 60.4 | – | – | – | – |
| NMS [22] | – | 63.2 | – | – | – | – |
| HAN [25] | – | 69.5 | 65.5 | 69.1 | – | – |
| LXMERT [47] | 59.8 | 60.0 | – | 72.4 | – | – |
| OSCAR [32] | – | 61.6 | – | 73.6 | – | – |
| UNITER-base [7] | – | – | – | 72.7 | – | – |
| UNITER-large [7] | – | – | – | **73.8** | – | – |
| **CFR (ours)** | **73.6** | **72.1** | **69.7** | 72.5 | **69.8** | **71.9** |

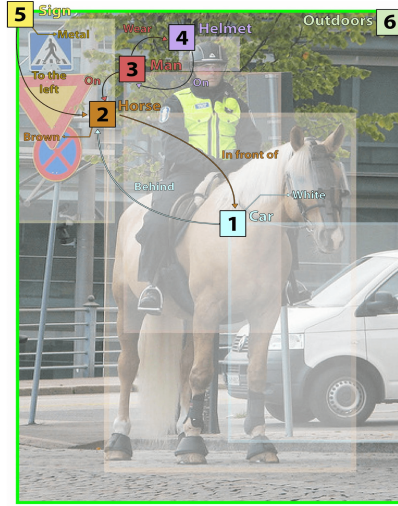Table 1. The accuracy of our method and other approaches on three VQA datasets.

and MCAN [59]; reasoning-based methods: Murel [4], ReGAT [30], MMN [6], NMS [22], and HAN [25]; and large-scale visual-language modeling: LXMERT [47], OSCAR [32], and UNITER [7].

**Evaluation Metrics.** As the standard practice, we use the accuracy metric ($Acc$) [3] to evaluate the free-form opened ended dataset (GQA and VQA 2.0), and $Acc$-$MC$ [65] to evaluate the multiple-choice dataset (Visual7W).
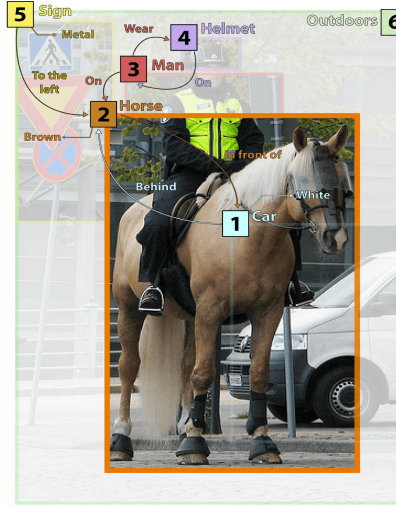
### 4.2. Module Contribution

### 4.3. Results

Table 1 summarizes our results compared with different recent methods in the VQA task. In the GQA dataset, our proposed method outperforms the recent approach HAN [25] on the test-dev set by $+2.6\%$. Regarding the multiple-choice Visual7W dataset, our method outperforms the work CTI [11] by $2.8\%$ in the validation set and $2.6\%$ in the test set, respectively. The results show that our CFR can deal with compositional reasoning questions through the selected information from both coarse-grained learning and fine-grained learning. It is worth noting that our CFR achieves new state-of-the-art results in GQA and Visual7W datasets.

(a)      (b)      (c)

(d)      (e)      (f)

Figure 3. Visualization of the explicit contribution of RoIs and predicates in both input image and question. The ✓ and ✗ symbols indicate the correct and the wrong answers, respectively. The arrow indicates the attribute or relation from the attribute classification or relation classification step in our Image Embedding module.

It is more challenging for our method to improve the result in the VQA2.0 dataset. While our CFR still outperforms the recent reasoning work ReGAT [30] by 2.5% and 2.2%, UNITER-large [7] achieves 1.3% higher than our CFR in the test-dev set. We note that the VQA2.0 dataset has fairly fewer compositional reasoning questions comparing with the GQA dataset [21]. Thus, it limits the effectiveness of methods that focus on reasoning the question and images, including our CFR. Our method also uses simple modules to extract image and question features,

| Methods | | Language Modality | | | Vision Modality | | | Acc (%) |
|---|---|---|---|---|---|---|---|---|
| | | *Question Features* | *Predicates* | *Filtered Info* | *Image Features* | *Predicates* | *Filtered Info* | |
| **Multimodal Learning** | *Coars grained* | ✓ | | | ✓ | | | 62.6 |
| | *Fine grained* | ✓ | ✓ | | ✓ | ✓ | | 67.2 (+4.6) |
| | | ✓ | | ✓ | ✓ | | ✓ | 69.5 (+6.9) |
| **Semantic Reasoning** | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 73.6 (+11.0) |

Table 2. The contribution of each module in our CFR framework.

which may not be robust enough comparing with features extracted from complicated modules such as large-scale visual-language models [32] [7].

To evaluate the contribution of each module in our framework, we conduct the following experiment: Given different level of information of language and vision modality (features, predicates, and filtered information of the image/question), we gradually choose different pairs of vision and language modality as the input to predict the answer. The experiment is conducted using the GQA dataset.

Table 2 shows the contribution of each module when different inputs are used. By using only the question and image feature (coarse-grained learning), our framework only achieves 62.6% accuracy. When we combine the question and image features with their corresponding predicates (fine-grained learning), the accuracy increase to 67.2%. This result indicates the effectiveness of predicates. By applying the filtered information of both question and image, the performance of fine-grained learning increases to 69.5%. This result shows that by reducing the negative influence of noisy information, the prediction accuracy can be improved. To effectively leverage all coarse-grained and fine-grained information, the Semantic Module is integrated into the framework and achieves 73.6% accuracy. This result validates the potency of Semantic Reasoning in selecting information for answering the complicated question. Overall, our introduced framework outperforms the baseline coarse-grained learning method by a large margin, i.e., +11.0% accuracy.

### 4.4. Visualization

Figure 2 illustrates the comparison between using Coarse-grained learning, Fine-grained learning, and Semantic Reasoning when we visualize the confidence score of the top 5 output answers. From this figure, we notice that if the Coarse-grained or Fine-grained learning are used separately, the output answer may not be correct, and there is usually an ambiguity in the top two predicted answers. However, when we apply our whole Coarse-to-Fine Reasoning framework, the network predicts both answers correctly, and also there is no ambiguity between the top pre-

dicted answer and the second predicted answer. These results show that our Coarse-to-Fine Reasoning framework successfully encodes both the features and predicates from the image and question in a coarse-to-fine manner, hence consequently improves the prediction results.

Figure 3 illustrates the explicit contribution of RoIs and predicates in both input image and question when our framework answers different compositional questions. Note that the transparency level of each RoI/word indicates the importance of each information. The RoIs and predicates with no opacity are crucial instances for answering the corresponding question. The visualizations in samples $(a, b, c, d, e)$ indicate the effectiveness of our CFR framework in reasoning the correct answers from the inference process. The sample in $(f)$ demonstrates the case when our CFR predicts the wrong answer. The incorrect prediction may come from the limitation of extractors, i.e., the extracted features are not robust enough (e.g., "cap" and "helmet" in our false example). Figure 3 also shows that our CRF framework not only can increase the accuracy of the VQA task but also provides an explainable way to understand the prediction results.

## 5. Conclusion

We have introduced a new simple, yet effective Coarse-to-Fine Reasoning (CFR) framework for the VQA task. Our CRF framework first extracts the features and predicates of both question and image. Then we propose a new reasoning module to map the key information in the question to the visual clues in the image in a coarse-to-fine manner. The intensive experiments on GQA, VQA2.0, and Visual7W datasets show that our framework achieves competitive results comparing with recent approaches. Our source code and trained models will be released for reproducibility and further study.

## References

[1] Saeed Amizadeh, Hamid Palangi, Oleksandr Polozov, Yichen Huang, and Kazuhito Koishida. Neuro-symbolic vi-

sual reasoning: Disentangling" visual" from" reasoning". In *ICML*, 2020. 2

[2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and VQA. In *CVPR*, 2018. 2

[3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *ICCV*, 2015. 5

[4] Remi Cadene, Hedi Ben-Younes, Matthieu Cord, and Nicolas Thome. Murel: Multimodal relational reasoning for visual question answering. In *CVPR*, 2019. 1, 2, 5

[5] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In *CVPR*, 2019. 1, 2

[6] Wenhu Chen, Zhe Gan, Linjie Li, Yu Cheng, William Wang, and Jingjing Liu. Meta module network for compositional visual reasoning. In *WACV*, 2021. 2, 5

[7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020. 2, 5, 6, 7

[8] Wenqiang Chi, Giulio Dagnino, Trevor MY Kwok, Anh Nguyen, Dennis Kundrat, Mohamed EMK Abdelaziz, Celia Riga, Colin Bicknell, and Guang-Zhong Yang. Collaborative robot-assisted endovascular catheterization with generative adversarial imitation learning. In *ICRA*, 2020. 2

[9] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*, 2014. 2

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 2

[11] Tuong Do, Thanh-Toan Do, Huy Tran, Erman Tjiputra, and Quang D Tran. Compact trilinear interaction for visual question answering. In *ICCV*, 2019. 2, 5

[12] Tuong Do, Binh X Nguyen, Erman Tjiputra, Minh Tran, Quang D Tran, and Anh Nguyen. Multiple meta-model quantifying for medical visual question answering. In *MICCAI*, 2021. 2

[13] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning approach for object affordance detection. In *ICRA*, 2018. 2

[14] Difei Gao, Ke Li, Ruiping Wang, Shiguang Shan, and Xilin Chen. Multi-modal graph neural network for joint reasoning on vision and scene text. In *CVPR*, 2020. 1, 2

[15] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *CVPR*, 2019. 1, 2, 5

[16] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *CVPR*, 2018. 1

[17] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 5

[18] Jiuxiang Gu, Handong Zhao, Zhe Lin, Sheng Li, Jianfei Cai, and Mingyang Ling. Scene graph generation with external knowledge and image reconstruction. In *CVPR*, 2019. 1, 2

[19] Xin Hong, Yanyan Lan, Liang Pang, Jiafeng Guo, and Xueqi Cheng. Transformation driven visual reasoning. In *CVPR*, 2021. 2

[20] Hexiang Hu, Wei-Lun Chao, and Fei Sha. Learning answer embeddings for visual question answering. In *CVPR*, 2018. 5

[21] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 5, 6

[22] Drew A. Hudson and Christopher D. Manning. Learning by abstraction: The neural state machine. In *NIPS*, 2019. 1, 5

[23] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. In defense of grid features for visual question answering. In *CVPR*, 2020. 2

[24] Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Pythia v0.1: the winning entry to the vqa challenge 2018. *CoRR*, 2018. 5

[25] Eun-Sol Kim, Woo Young Kang, Kyoung-Woon On, Yu-Jung Heo, and Byoung-Tak Zhang. Hypergraph attention networks for multimodal learning. In *CVPR*, 2020. 5

[26] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In *NIPS*, 2018. 1, 2, 3, 5

[27] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, pages 32–73, 2016. 2, 5

[28] Thao Minh Le, Vuong Le, Svetha Venkatesh, and Truyen Tran. Dynamic language binding in relational visual reasoning. In *IJCAI*, 2020. 2

[29] Guohao Li, Xin Wang, and Wenwu Zhu. Perceptual visual reasoning with knowledge propagation. In *Proceedings of the 27th ACM International Conference on Multimedia*, 2019. 1

[30] Linjie Li, Zhe Gan, Yu Cheng, and Jingjing Liu. Relation-aware graph attention network for visual question answering. In *ICCV*, 2019. 2, 5, 6

[31] Linjie Li, Jie Lei, Zhe Gan, and Jingjing Liu. Adversarial vqa: A new benchmark for evaluating the robustness of vqa models. *arXiv:2106.00245*, 2021. 2

[32] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 2, 5, 7

[33] Xiaodan Liang, Lisa Lee, and Eric P Xing. Deep variation-structured reinforcement learning for visual relationship and attribute detection. In *CVPR*, 2017. 1

[34] Edward Loper and Steven Bird. Nltk: the natural language toolkit. *arXiv preprint cs/0205028*, 2002. 2

[35] Pan Lu, Lei Ji, Wei Zhang, Nan Duan, Ming Zhou, and Jianyong Wang. R-vqa: learning visual relation facts with semantic attention for visual question answering. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018. 1

[36] Man Luo, Shailaja Keyur Sampat, Riley Tallman, Yankai Zeng, Manuha Vancha, Akarshan Sajja, and Chitta Baral. 'just because you are right, doesn't mean i am wrong': Overcoming a bottleneck in development and evaluation of open-ended vqa tasks. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021. 2

[37] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *ICLR*, 2019. 2

[38] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *CVPR*, 2019. 1

[39] David Mascharka, Philip Tran, Ryan Soklaski, and Arjun Majumdar. Transparency by design: Closing the gap between performance and interpretability in visual reasoning. In *CVPR*, 2018. 2

[40] Medhini Narasimhan, Svetlana Lazebnik, and Alexander Schwing. Out of the box: Reasoning with graph convolution nets for factual visual question answering. In *NIPS*, 2018. 1

[41] Anh Nguyen, Thanh-Toan Do, Ian Reid, Darwin G Caldwell, and Nikos G Tsagarakis. V2cnet: A deep learning framework to translate videos to commands for robotic manipulation. *arXiv:1903.10869*, 2019. 2

[42] Anh Nguyen, Ngoc Nguyen, Kim Tran, Erman Tjiputra, and Quang D Tran. Autonomous navigation in complex environments with deep multimodal fusion network. In *IROS*, 2020. 1

[43] Anh Nguyen, Quang D Tran, Thanh-Toan Do, Ian Reid, Darwin G Caldwell, and Nikos G Tsagarakis. Object captioning and retrieval with natural language. In *CVPRW*, 2019. 2

[44] Binh X Nguyen, Binh D Nguyen, Tuong Do, Erman Tjiputra, Quang D Tran, and Anh Nguyen. Graph-based person signature for person re-identifications. In *CVPRW*, 2021. 1

[45] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014. 2, 5

[46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 2

[47] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *EMNLP*, 2019. 1, 2, 5

[48] Damien Teney, Ehsan Abbasnejad, and Anton van den Hengel. Unshuffling data for improved generalization in visual question answering. In *ICCV*, 2021. 2

[49] Damien Teney, Peter Anderson, Xiaodong He, and Anton van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. In *CVPR*, 2018. 2, 5

[50] Damien Teney, Lingqiao Liu, and Anton van den Hengel. Graph-structured representations for visual question answering. In *CVPR*, 2017. 1, 2

[51] Minh Tran, Tuong Do, Binh X Nguyen, Erman Tjiputra, Quang D Tran, and Anh Nguyen. Light-weight deformable registration usingadversarial learning with distilling knowledge. *arXiv:2110.01293*, 2021. 2

[52] Aisha Urooj, Hilde Kuehne, Kevin Duarte, Chuang Gan, Niels Lobo, and Mubarak Shah. Found a reason for me? weakly-supervised grounded visual question answering using capsules. In *CVPR*, 2021. 2

[53] Peng Wang, Qi Wu, Chunhua Shen, Anthony R. Dick, and Anton van den Hengel. Fvqa: Fact-based visual question answering. *TPAMI*, pages 2413–2427, 2017. 1

[54] Zhe Wang, Xiaoyi Liu, Limin Wang, Yu Qiao, Xiaohui Xie, and Charless Fowlkes. Structured triplet learning with pos-tag guided attention for visual question answering. In *WACV*, 2018. 5

[55] Jialin Wu and Raymond J Mooney. Self-critical reasoning for robust visual question answering. In *NIPS*, 2019. 2

[56] Xiaofeng Yang, Guosheng Lin, Fengmao Lv, and Fayao Liu. Trrnet: Tiered relation reasoning for compositional visual question answering. In *ECCV*, 2020. 1

[57] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alexander J. Smola. Stacked attention networks for image question answering. In *CVPR*, 2016. 2, 5

[58] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Josh Tenenbaum. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. In *NIPS*, 2018. 2

[59] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *CVPR*, 2019. 2, 5

[60] Weifeng Zhang, Jing Yu, Hua Hu, Haiyang Hu, and Zengchang Qin. Multimodal feature fusion by relational reasoning and attention for visual question answering. *Information Fusion*, 2020. 1

[61] Weifeng Zhang, Jing Yu, Wenhong Zhao, and Chuan Ran. Dmrfnet: Deep multimodal reasoning and fusion for visual question answering and explanation generation. *Information Fusion*, 2021. 2

[62] Chen Zheng, Quan Guo, and Parisa Kordjamshidi. Cross-modality relevance for reasoning on language and vision. In *ACL*, 2020. 2

[63] Wenbo Zheng, Lan Yan, Chao Gou, and Fei-Yue Wang. Webly supervised knowledge embedding model for visual reasoning. In *CVPR*, 2020. 2

[64] Wenfeng Zheng, Lirong Yin, Xiaobing Chen, Zhiyang Ma, Shan Liu, and Bo Yang. Knowledge base graph embedding module design for visual question answering model. *Pattern Recognition*, 2021. 2

[65] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7W: Grounded question answering in images. In *CVPR*, 2016. 2, 5