# Question Guided Modular Routing Networks for Visual Question Answering

**Yanze Wu[1], Qiang Sun[1], Jianqi Ma[1], Bin Li[1], Yanwei Fu[1], Yao Peng[2], Xiangyang Xue[1]**

[1]Fudan University, [2]Qiniu Inc

{wuyanze123,sunqiang85}@gmail.com, {majq16,libin,yanweifu}@fudan.edu.cn, pengyao@qiniu.com, xyxue@fudan.edu.cn

## Abstract

This paper studies the task of Visual Question Answering (VQA), which is topical in Multimedia community recently. Particularly, we explore two critical research problems existed in VQA: (1) efficiently fusing the visual and textual modalities; (2) enabling the visual reasoning ability of VQA models in answering complex questions. To address these challenging problems, a novel Question Guided Modular Routing Networks (QGMRN) has been proposed in this paper. Particularly, The QGMRN is composed of visual, textual and routing network. The visual and textual network serve as the backbones for the generic feature extractors of visual and textual modalities. QGMRN can fuse the visual and textual modalities at multiple semantic levels. Typically, the visual reasoning is facilitated by the routing network in a discrete and stochastic way by using Gumbel-Softmax trick for module selection. When the input reaches a certain modular layer, routing network newly proposed in this paper, dynamically selects a portion of modules from that layer to process the input depending on the question features generated by the textual network. It can also learn to reason by routing between the generic modules without additional supervision information or expert knowledge. Benefiting from the dynamic routing mechanism, QGMRN can outperform the previous classical VQA methods by a large margin and achieve the competitive results against the state-of-the-art methods. Furthermore, attention mechanism is integrated into our QGMRN model and thus can further boost the model performance. Empirically, extensive experiments on the CLEVR and CLEVR-Humans datasets validate the effectiveness of our proposed model, and the state-of-the-art performance has been achieved.

## 1 Introduction

Visual Question Answering (VQA) (Antol et al. 2015) is a task that, given an image and question pair, the model can reason and answer the question by visual information. It is a popular topic in recent years that has connected the computer vision and natural language processing (NLP). VQA faces two major challenges: 1) How to better fuse the visual and textual modalities; 2) How to make the VQA model have the reasoning ability to answer more complex questions.

For challenge (1), in order to fuse the visual and textual modalities, the most common VQA paradigm is to first ex-

tract the visual features from a modern CNN (*e.g.,* VGG (Simonyan and Zisserman 2014) or ResNet (He et al. 2016)) and textual features from a RNN (*e.g.,* LSTM (Hochreiter and Schmidhuber 1997) or GRU (Cho et al. 2014)) separately, and then fuse them into a common latent space (Kazemi and Elqursh 2017; Anderson et al. 2018; Malinowski, Rohrbach, and Fritz 2015; Fukui et al. 2016; Wu et al. 2018a). Feature fusion is explored from simple concatenation operation (Zhou et al. 2015) to advanced joint embedding techniques, such as MLB (Kim et al. 2016), MCB (Fukui et al. 2016) and MUTAN (Ben-Younes et al. 2017), and attention mechanisms such as one-hop attention (Kazemi and Elqursh 2017), multi-hop attention (Yang et al. 2016) and bottom-up attention (Anderson et al. 2018). However, no matter how powerful these fusion methods are, they can only fuse two modalities at a high semantic level, rather than multiple semantic levels, because the fused features (*i.e.,* the extracted visual and textual features) are at a high semantic level.

For challenge (2), to endow the VQA model with reasoning ability is difficult because although CNN (*e.g.,* ResNet) is very powerful, it does not support reasoning natively. To make the VQA model support reasoning, one important line of work (Andreas et al. 2016; Hu et al. 2017; Johnson et al. 2017b; Mascharka et al. 2018; Mao et al. 2019) builds upon the modular networks. By leveraging the compositional structure of neural language questions, these methods first decompose the question into modular subquestions, then pick out the modular networks that are designed for solving these subquestions, and finally answer the question by executing these modular networks on the image according to the layout learned from the question. These explicit methods have strong interpretability as the layout learned from the question reflects the answering and reasoning process. However, analyzing the subquestion set and designing corresponding modules need expert knowledge, most of these methods even require extra supervision information during their training process.

In this paper, we show that even without expert knowledge and extra supervision information, we can still make the model support visual reasoning. Specifically, we propose a novel model called Question Guided Modular Routing Net-

works (QGMRN) that consists of three sub-networks: *visual network*, *textual network* and *routing network*. The visual network is based on a generic modular network that each layer of the network is composed of some generic modules, as the module granularity changes, different modular networks can be spawned. When the input reaches a certain layer, the visual network dynamically selects a portion of modules from that layer to process the input according to the binary gates generated by the routing network. We name the collection of all binary gates as a routing path. Therefore, the routing network is responsible for receiving the question features generated by the textual network and mapping them to a discrete routing path. To encourage the model to explore more possible routing paths, the model generates the routing path in a discrete and stochastic way by adopting Gumbel-Softmax tricks. We claim that our proposed method can address both two before-mentioned challenges to some extent. First, the visual and textual modalities can be fully fused at multiple semantic levels because the routing network controls every layer of the visual network. Second, the routing mechanism makes the model have the ability for visual reasoning by executing different portion of modules. Especially we claim that the routing path represents the process of compositional reasoning and acts similar to the aforementioned module layout.

To sum up, the contributions of this paper can be summarized as: (1) We for the first time successfully extend the line of work on routing models to the domain of VQA and multi-modal embedding, and the results are promising. (2) Our model can fuse the visual and textual modalities at multiple semantic levels, and is equipped with good reasoning ability. (3) Current modular network based VQA methods have great interpretability, but they also have the disadvantages of requiring expert knowledge or extra supervision information. Without having those disadvantages, we combine the routing mechanism with the modular network to obtain the modular routing network. The novel modular routing network can be applied to many off-the-shelf CNN models. (4) The proposed model has achieved state-of-the-art performances on the challenging CLEVR (Johnson et al. 2017a) and CLEVR-Humans (Johnson et al. 2017b) datasets.

## 2 Related Work

### 2.1 Fusion Strategy of VQA

Usually, the first step of most VQA methods is to extract high-level visual features from a modern CNN and textual features from an RNN separately. In order to combine the visual and textual modalities to produce the answer, many methods have been proposed to fuse the extracted visual and textual features. Multimodal Low-rank Bilinear pooling (MLB) (Kim et al. 2016) provides an efficient method to approximate the full bilinear pooling by forcing the weight matrix to be low-rank. Multimodal Compact Bilinear pooling (MCB) (Fukui et al. 2016) randomly projects the visual and textual features into a higher dimensional space, then processes them in Fast Fourier Transform space. Multimodal Tucker Fusion (MUTAN) (Ben-Younes et al. 2017) proposes a general fusion method based on Tucker decom-

position, which covers MLB and MCB. Yang *et al.* (Yang et al. 2016) propose a multi-hop spatial attention to fuse the visual and textual features so that the image regions related to the question will be focused. Anderson *et al.* (Anderson et al. 2018) propose a bottom-up model that combines the attention mechanism with object-level visual features, so that objects related to the question will be focused. However, with only the high-level multi-modal features fused, these fusion methods do not fully utilized the multi-level interactions between two modalities.

Another line of work proposes to fuse two modalities by using the question to predict the parameters of the visual network. Due to the large number of parameters in the visual network, only a small portion of parameters can be learned from the question feature. For example, Gao *et al.* (Gao et al. 2018) propose the Question-guided Hybrid Convolution (QGHC) based on group convolution, which consists of question-dependent kernels and question-independent kernels; the parameters of batch-norm (Ioffe and Szegedy 2015) layers in MODERN (De Vries et al. 2017) are predicted by the question; FiLM (Perez et al. 2018) implements the condition through a general feature-wise transformation. Although this kind of methods can fuse the visual and textual modalities in a multi-level and fine-grained way, these methods still have limitations. If too many parameters are learned from the question, that will make the model difficult to train, but if only a few parameters are learned from the question, that will constrain the model's learning capacity. Another concern is about the flexibility. In particular, QGHC can only be applied in the ResNext (Xie et al. 2017) architecture, MODERN can only be applied in CNN with the batch-norm layers.

Our proposed model can fuse the two modalities at multiple semantic levels, and more importantly, our model is very flexible. When we adjust the granularity of the module to filter level, MODERN and FiLM are related to our model, yet with fundamental differences. First, they modulate the visual network at a feature level, but we modulate the visual network at a module level. Second, we route the modules in a discrete and stochastic way, which means that our approach can explore more routing paths and is a step towards discrete reasoning.

### 2.2 Visual Reasoning in VQA

Judging whether a VQA model has reasoning ability is usually based on whether the model can well solve the CLEVR dataset, since the questions in CLEVR dataset are quite complicated that require reasoning to answer. There have been several prominent models that well solve the CLEVR dataset. Except for RN (Santoro et al. 2017) that supports reasoning explicitly by utilizing pairwise comparisons, most methods are based on the idea of compositional reasoning, implicitly (Perez et al. 2018; Hudson and Manning 2018) or explicitly (Andreas et al. 2016; Hu et al. 2017; Johnson et al. 2017b; Mascharka et al. 2018; Mao et al. 2019). The representative method of explicit compositional reasoning is Neural Module Networks (NMN) (Andreas et al. 2016) which is dynamically instantiated from a collection of reusable modules based on the compositional struc-

ture of the question. Although the function of each module is learned from training, the question parser and the mapping rules from the parsing tree node to the module must be pre-defined. The performance of the NMN model heavily relies on the quality of the question parser chosen. Further, Hu *et al.* (Hu et al. 2017) proposes an End-to-End Module Network which predicts the module layout by an LSTM instead of an external question parser. Johnson *et al.* (Johnson et al. 2017b) proposes a model combining with both program generator and execute engine based on neural module network. TbD-net (Mascharka et al. 2018) combines neural modules and attention mechanism to achieve state-of-the-art accuracy on the CLEVR dataset. Different from above methods, NS-CL (Mao et al. 2019) can achieve state-of-the-art performance even without extra supervision information.

Although the NMN family has good interpretability, they usually require expert knowledge for designing the model and extra supervision information for training. Our method shares the concept of the module with the NMN family, however, the modules from our method are not dedicated-designed and no extra supervision information is required for training, we also combine the modular network with the routing mechanism. Compared with the explicit reasoning models like FiLM (Perez et al. 2018) and MAC (Hudson and Manning 2018), our method uses a completely different routing mechanism to support reasoning.

## 2.3 Routing Models

The proposed modular routing network is related to conditional computation (CC) (Bengio, Léonard, and Courville 2013; Bengio et al. 2016) and mixture-of-experts (MoE) (Jordan and Jacobs 1994; Jacobs et al. 1991; Eigen, Ranzato, and Sutskever 2013; Shazeer et al. 2017). CC refers to the dynamic execution of a part of the model based on the input. Bengio *et al.* (Bengio, Léonard, and Courville 2013) first propose the concept of CC and introduce four approaches to propagate gradients through stochastic neurons to support CC. Bengio *et al.* (Bengio et al. 2016) propose by using reinforcement learning as a tool to optimize CC, they also propose a regularization mechanism to encourage sparsity and diversity. MoE related research can be traced back to the 1990s (Jordan and Jacobs 1994; Jacobs et al. 1991), and usually, an expert is a whole model. Eigen *et al.* (Eigen, Ranzato, and Sutskever 2013) propose a deep learning model that is made up of two MoEs, Shazeer *et al.* (Shazeer et al. 2017) introduce the sparsely-gated MoE by using a noisy gating method. Kirsch *et al.* (Kirsch, Kunze, and Barber 2018) use generalized Viterbi EM to enable training without artificial regularization. Routing model is equivalent to doing CC on MoE. Researchers have tried to apply routing models to multi-task learning (Rosenbaum, Klinger, and Riemer 2018) and classification problems (Wu et al. 2018b; Veit and Belongie 2018). However, these methods are either adopting small-scale models, experimenting on toy-dataset, or achieving little improvement on the basis of static models.

As Ramachandran *et al.* said in their paper (Ramachandran and Le 2019), although routing models are a promising direction of research, there must be successful applications of routing models where static models struggle. In this pa-

per, we find that VQA (especially those who need strong reasoning ability) is one of the successful applications, the static models (*i.e.,* standard VQA methods) perform poorly on the VQA datasets that require reasoning, but the routing models outperform them by a large margin.

## 3 METHODOLOGY

Our model contains three sub-networks, namely the *visual network*, *textual network* and *routing network* as shown in Figure 1. The textual network takes a question and generates question features; the routing network converts the question features to a specific routing path; the visual network takes the raw image and dynamically executes modules according to the routing path. Finally, the extracted image features are sent to the classifier to predict the answer. The whole model is differentiable with respect to the model parameters and can be trained end-to-end.

### 3.1 Textual Network

Given the question $\mathbf{Q} = [\mathbf{w}_1, \mathbf{w}_2, \cdots, \mathbf{w}_n]$ where $\mathbf{w}_i$ is the one-hot representation of $i$th word and $n$ is the length of the question, we first use a lookup layer to embed $\mathbf{Q}$ into $\mathbf{E}_q = [\mathbf{e}_1, \mathbf{e}_2, \cdots, \mathbf{e}_n]$ where $\mathbf{e}_i \in \mathbb{R}^d$ , $d$ is the embedding size. Then we feed $\mathbf{E}_q$ into Gated Recurrent Units (GRU) and use the final hidden state of GRU as the question features:

$$\mathbf{q} = \text{GRU}(\mathbf{E}_q) \qquad (1)$$

where $\mathbf{q} \in \mathbb{R}^h$ and $h$ is the hidden size of GRU units.

### 3.2 Visual Network

To better illustrate our approach, we will first introduce the generic Modular Architecture and Modular Routing Architecture. Modular Architecture defines the granularity and organization structure of the modules. Module Routing Architecture defines how the routing path controls over the Modular Architecture. Then we will introduce two forms of implementations of the module routing architecture: **FRN**, which is based on ResNet; **BRN** (**B**ranch **R**outing **N**etwork), which is based on ResNext (Xie et al. 2017).

**Modular Architecture.** We assume that a generic modular architecture is composed of $L$ modular layers, each modular layer is composed of $M$ modules[1], each module $f_{l,m}(x; \theta_{l,m})$ for $l \in \{1, \cdots, L\}$ and $m \in \{1, \cdots, M\}$ is a function module that takes the input $x$ and generates an output tensor, $\theta_{l,m}$ is the learnable parameter of $f_{l,m}$. All modules in the $l$th layer share the same input, and the output of the $l$th layer $y_l$ is:

$$y_l = \phi([f_{l,1}(y_{l-1}; \theta_{l,1}), \cdots, f_{l,M}(y_{l-1}; \theta_{l,M})]) + y_{l-1} \quad (2)$$

where the composite function $\phi(.)$ can either be concatenation or summation.

---

[1]For convenience, we assume that the number of modules per modular layer is the same, but in fact they can be different.
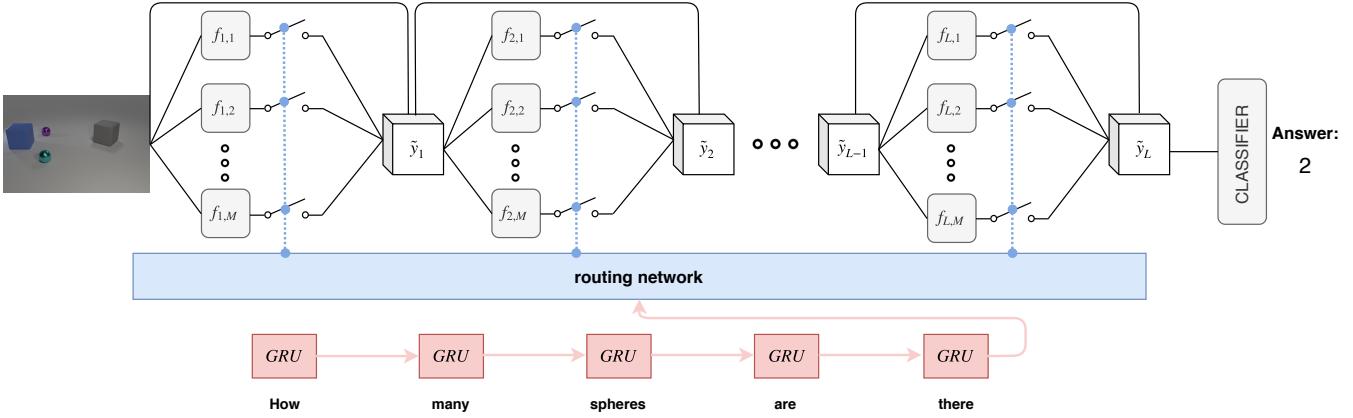
Figure 1: Architecture of the proposed model. The model consists three sub-networks: a) Visual Network: a composite of $L \times M$ modules which are represented as $f_{1..L,1..M}$, we can see that each module is followed by a "switch" which is controlled by the Routing Network, we call the states of the "switches" as a routing path. b) Textual Network: a GRU which takes a question and generates question features. c) Routing Network: which takes the question features sent from the textual network to generate a routing path.

**Modular Routing Architecture**    First, we denote the routing path as

$$\mathbf{P} = \text{RNET}(\mathbf{q}) \tag{3}$$

where RNET is the routing network that generates a routing path $\mathbf{P} \in \{0,1\}^{L \times M}$ condtioned on the question features $\mathbf{q} \in \mathbb{R}^h$. $\mathbf{P}_{l,m}$ for $l \in \{1, \cdots, L\}$ and $m \in \{1, \cdots, M\}$ is a binary gate that controls whether or not to execute the $m$th module of the $l$th layer. The output of the $l$th layer $\tilde{y}_l$ of modular routing architecture is now changed to:

$$\tilde{y}_l = \phi([\tilde{f}_{l,1}(\tilde{y}_{l-1}; \theta_{l,1}), \cdots, \tilde{f}_{l,M}(\tilde{y}_{l-1}; \theta_{l,M})]) + \tilde{y}_{l-1} \tag{4}$$

where $\tilde{f}_{l,m}(\tilde{y}_{l-1}; \theta_{l,m}) = \mathbf{P}_{l,m} \cdot f_{l,m}(\tilde{y}_{l-1}; \theta_{l,m})$. Due to the discrete nature of $\mathbf{P}$, the gradient backpropagation algorithm cannot be applied here, we will use some tricks to make the whole model differentiable and the details will be discussed in section 3.3.

**FRN and BRN**    As we claimed before, with the introduced modular routing architecture, different models can be realized by changing the granularity of the module. Here we introduce two kinds of special cases that are easy to implement.

**FRN:** when the routed module is the filter in a convolutional layer, we call it FRN. Note that the FRN can be plugged into any modern CNN, we applied it to the current popular ResNet for convenience. More specifically, we route the last convolutional layer of each *residual block*.

**BRN:** when the routed module is the branch of a multi-branch network, we call it BRN. The BRN can be plugged into any multi-branch convolutional networks. In this paper, we test the BRN under the ResNext model.

### 3.3   Routing Network

As we claimed before, RNET aims to compute the binary gates conditioned on the question features. First, we temporarily ignore the fact that the routing path is discrete for simplicity, All we need is a fully connected layer fc : $\mathbb{R}^h \rightarrow \mathbb{R}^{L \times M}$ to map the question features to a *real-valued* routing path $\tilde{\mathbf{P}} = \text{fc}(\mathbf{q})$. Now take into account the nature of the discrete, a naive attempt would be thresholding $\tilde{\mathbf{P}}$ into 1s and 0s, but unfortunately, it is not differentiable and the backpropagation algorithm cannot be applied. Also note that the threshold function is deterministic, in order to explore more possible paths, the generation of routing path will better to be stochastic. Based on the above considerations, *i.e.,* **discrete** and **stochastic**, we let the routing path $\mathbf{P}$ follows a $L \times M$-dimensional Bernoulli distribution whose parameter $\mathbf{s} \in [0,1]^{L \times M}$ is generated from the routing network. To optimize the whole model, we employ a reparameterization trick called *Concrete Distribution* (Maddison, Mnih, and Teh 2017) or *Gumbel- Softmax* (Jang, Gu, and Poole 2017) in this paper.

In order to elicit this method, we first review the Gumbel-Max trick (Yellott Jr 1977) which provides a way to sample $\mathbf{z}$ from a categorical distribution with class probability of $\pi_1, \pi_2, \cdots, \pi_n$ as follows:

$$\mathbf{z} = \text{one-hot}(\arg \max_i [\mathbf{g}_i + \log \pi_i]) \tag{5}$$

where $\mathbf{g}_1, ..., \mathbf{g}_n$ are i.i.d samples from Gumbel distribution [2] and $P(\mathbf{z}_k = 1) = \pi_k$. However the argmax operation is still not differentiable, so the softmax function with temperature $\tau$ is introduced here to approximate the argmax function:

$$\tilde{\mathbf{z}} = \text{softmax}((\mathbf{g} + \log(\pi))/\tau) \tag{6}$$

As $\tau \rightarrow 0$, the softmax function is smoothly approaching the argmax function.

In the above we have showed how to sample from categorical distribution $\text{Cat}(\pi_1, \pi_2, \cdots, \pi_n)$ using the Gumbel-

---

[2]To sample from Gumbel distribution, first draw a sample from Uniform distribution $u \sim \text{Uniform}(0,1)$, computing random variable g as $g = -\log(-\log(u))$, then $g \sim \text{Gumbel}(0,1)$

Softmax trick, now we discuss the "binary" case of our problem, that means we need to sample from $\text{Cat}(\pi_1, \pi_2)$ where $\pi_2 = 1 - \pi_1$. Note that this distribution is equivalent to Bernoulli distribution $\text{Bern}(\pi_1)$, we use $\rho$ to substitute $\pi_1$ to distinguish from the previous symbol definition, then we apply the Gumbel-Max trick again to sample Bernoulli variable b from $\text{Bern}(\rho)$ as

$$\text{b} = [\text{g}_1 + \log(\rho) > \text{g}_2 + \log(1 - \rho)] \quad (7)$$
$$= [\text{g}_1 - \text{g}_2 + \log(\rho/(1 - \rho)) > 0] \quad (8)$$
$$= [\text{t} + \log(\rho/(1 - \rho)) > 0] \quad (9)$$

where $\text{t} \sim \text{Logistic}(0, 1)$ and the bracketed notation [statement] stands for 1 if statement is true, 0 otherwise. The derivation of Eq.(9) uses the fact that the difference of two Gumbels variables follows a Logistic distribution[3], *i.e.,* $\text{g}_1 - \text{g}_2 \sim \text{Logistic}(0, 1)$. Just as we use the softmax function with temperature to approximate the argmax function, here we use the sigmoid function with temperature to approximate the unit step function:

$$\tilde{\text{b}} = \text{sigmoid}((\text{t} + \log(\rho/(1 - \rho)))/\tau) \quad (10)$$

as $\tau \to 0$, the sigmoid function is smoothly approaching the unit step function. Readers who are interested in Gumbel-Softmax or Concrete Distribution are referred to (Maddison, Mnih, and Teh 2017; Jang, Gu, and Poole 2017) for more details.

Based on previous discussion, we can convert the *real-valued* $\tilde{\mathbf{P}}$ to binary value in a simple way. First, for each entry $\tilde{\mathbf{P}}_{l,m}$, we compute $\tilde{\mathbf{B}}_{l,m}$ as:

$$\tilde{\mathbf{B}}_{l,m} = \text{sigmoid}((\text{t}_{l,m} + \tilde{\mathbf{P}}_{l,m})/\tau) \quad (11)$$

where $\text{t}_{1...L,1...M}$ are i.i.d samples from Logistic distribution and $\tilde{\mathbf{P}}_{l,m} = \log(\frac{\mathbf{s}_{l,m}}{1 - \mathbf{s}_{l,m}})$. However $\tilde{\mathbf{B}}_{l,m}$ is still a continuous value as $\tau > 0$, here we use a Straight-Through (ST) method introduced in (Bengio, Léonard, and Courville 2013) to convert the continuous $\tilde{\mathbf{B}}_{l,m}$ to discrete $\mathbf{P}_{l,m}$, that is, during forward process, we use a threshold of $0.5$ to thresholding $\tilde{\mathbf{B}}_{l,m}$ to 0 and 1, but during the backward process, the gradient is normally passed to $\tilde{\mathbf{B}}_{l,m}$ just as the thresholding function is an identity function. We provide the pseudocode of the routing algorithm in appendix B.

## 3.4 QGMRN with Attention

Benefiting from the routing mechanism, the visual feature map $\tilde{y}_L$ extracted from the visual network already contains enough information and clues to answer the question, so we can directly send it to a global max pooling (GMP) layer and a simple MLP classifier to predict the answer.

And we can also use an additional attention mechanism to further fuse the multi-modal features. The attention method consists of two stages: (1) We use the spatial self-attention

mechanism to enhance the representation power and the spatial reasoning ability of the model; (2) We use a simple spatial attention mechanism to convert the feature map to the final feature vector.

In the first stage, we first concatenate the question features $\mathbf{q}$ with the visual features $\tilde{y}_L$ to get the fused multi-modal features $u \in \mathbb{R}^{h \times w \times c}$ (broadcasting is used here) where $h$, $w$, and $c$ denote the height, width and channel of the feature map. Then we reshape it to $hw \times c$, and feed it to an off-the-shelf Transformer (Vaswani et al. 2017) encoder layers for modeling the correlations among the spatial locations. After we get the output $\tilde{u} \in \mathbb{R}^{hw \times d_t}$ where $d_t$ is the dimension of the Transformer, we reshape it back to $h \times w \times d_t$.

In the second stage, rather than using GMP as aggregator to convert the feature map $\tilde{u}$ to the feature vector, we find the weighted summation as aggregator is better. Specifically, we send $\tilde{u}$ to a $1 \times 1$ convolution layer followed by a spatial softmax layer to get the attention/weight map $a \in \mathbb{R}^{h \times w \times 1}$, the final features $f$ send to the MLP classifier is the weighted sum of $\tilde{u}$ at all spatial locations:

$$f = \text{sum}(a \cdot \tilde{u}, \dim = (0, 1)) \quad (12)$$

where $f \in \mathbb{R}^{d_t}$.

## 3.5 Training Loss

In order to avoid model collapse and prevent certain modules from being always executed or always not executed, some sparsity and variance regularizations (Bengio et al. 2016; Veit and Belongie 2018; Shazeer et al. 2017) are introduced. In this paper, we do not require our model to be sparse. But we observe that adding an regularization technique called load-balancing loss $\mathcal{L}_{load}$, which was first introduced in (Shazeer et al. 2017) to prevent model collapse, will make the model converge faster. We define the load of a module as the number of samples in one training batch activate that module. So the load of the $m$th module in the $l$th layer is calculated as: $\mathbf{A}_{l,m} = \sum_{i=1}^{N} \mathbf{P}_{l,m}^{(i)}$, where $\mathbf{P}_{l,m}^{(i)}$ is the gate of the $m$th module in the $l$th layer for the $i$th instance in the mini-batch, $N$ is the batch size. Further, we define $\mathcal{L}_{load_l} = \text{CV}(\mathbf{A}_{l,:})^2$ as the square of the coefficient of variation (CV) for the module load at the $l$th layer, and $\mathcal{L}_{load} = \sum_{l=1}^{L} \mathcal{L}_{load_l}$ as the sum of $\mathcal{L}_{load_l}$ in all layers.

With the standard VQA loss $\mathcal{L}_{vqa}$, the full loss is:

$$\mathcal{L} = \mathcal{L}_{vqa} + \lambda \mathcal{L}_{load} \quad (13)$$

where $\lambda$ is the coefficient of the load-balancing loss.

# 4 EXPERIMENTS

## 4.1 Datasets

The proposed method is evaluated on two datasets: (1) CLEVR dataset (Johnson et al. 2017a) is proposed to study the ability of VQA systems to perform reasoning. Answering question about a CLEVR image requires various kinds of reasoning, which makes the standard VQA methods perform poorly on this dataset. The dataset contains 100K 3D-rendered images and about one million automatically-generated questions. Specifically, the question in the dataset can be divided into the following five types: `Count`: ask the

---

[3]To sample from Logistic distribution, computing random variable t as $\text{t} = \log(\text{u}) - \log(1 - \text{u})$ where $\text{u} \sim \text{Uniform}(0, 1)$, then $\text{t} \sim \text{Logistic}(0, 1)$

| Model | CLEVR Overall | Count | Exist | Compare Numbers | Query Attribute | Compare Attribute | Humans After FT |
|---|---|---|---|---|---|---|---|
| Human (Johnson et al. 2017b) | 92.6 | 86.7 | 96.6 | 86.5 | 95.0 | 96.0 | – |
| N2NMN* (Hu et al. 2017) | 83.7 | 68.5 | 85.7 | 84.9 | 90.0 | 88.7 | – |
| PG+EE (18K prog.)* (Johnson et al. 2017b) | 95.4 | 90.1 | 97.3 | 96.5 | 97.4 | 98.0 | 66.6 |
| PG+EE (700K prog.)* (Johnson et al. 2017b) | 96.9 | 92.7 | 97.1 | 98.7 | 98.1 | 98.9 | – |
| TbD* (Mascharka et al. 2018) | 99.1 | 97.6 | 99.2 | 99.4 | 99.5 | 99.6 | – |
| LSTM (Johnson et al. 2017b) | 46.8 | 41.7 | 61.1 | 69.8 | 36.8 | 51.8 | 36.5 |
| CNN+LSTM (Johnson et al. 2017b) | 52.3 | 43.7 | 65.2 | 67.1 | 49.3 | 53.0 | 43.2 |
| CNN+LSTM+SA (Johnson et al. 2017b) | 76.6 | 64.4 | 82.7 | 77.4 | 82.6 | 75.4 | 57.6 |
| FiLM (Perez et al. 2018) | 97.6 | 94.3 | 99.3 | 93.4 | 99.3 | 99.3 | 75.9 |
| QGHC (Gao et al. 2018) | 86.3 | 78.1 | 91.7 | 80.7 | 89.4 | 86.8 | – |
| Relation Network (Santoro et al. 2017) | 95.5 | 90.1 | 97.8 | 93.6 | 97.9 | 97.1 | – |
| MAC (Hudson and Manning 2018) | 98.9 | 97.2 | 99.5 | 99.4 | 99.3 | 99.5 | 81.5 |
| NS-CL (Mao et al. 2019) | 98.9 | 98.2 | 98.9 | 99.0 | 99.3 | 99.1 | – |
| BRN (c=8) (ours) | 86.8 | 79.5 | 93.4 | 75.1 | 90.0 | 91.1 | – |
| BRN (c=16) (ours) | 94.7 | 90.8 | 98.1 | 82.8 | 97.3 | 98.4 | – |
| BRN (c=32) (ours) | 97.9 | 95.4 | 98.9 | 92.9 | 99.4 | 99.7 | 77.9 |
| BRN (c=32)+Attention (ours) | 98.7 | 97.4 | 99.6 | 95.1 | 99.7 | 99.7 | 79.3 |
| FRN (ours) | 98.2 | 96.6 | 99.9 | 94.1 | 98.7 | 99.9 | 79.9 |
| FRN+Attention (ours) | 98.9 | 97.7 | 99.9 | 94.5 | 99.8 | 99.9 | 81.8 |

Table 1: Comparison of accuracy on CLEVR and CLEVR-Humans datasets with previous methods. (*) denotes use of extra supervisory information through program labels. And c is the abbreviation for cardinality, and represents the number of modules per modular layer. The highest and second highest accuracies are shown in different background colors. However, we did not highlight TbD since it uses the extra supervisory information.

number of certain objects; Exist: ask whether a certain object is present; Compare Numbers: ask which of two object sets is larger; Query Attribute: query a attribute of particular object; Compare Attribute: ask whether two particular objects have same value on some attribute. (2) CLEVR-Humans dataset (Johnson et al. 2017b) contains human-posed questions on CLEVR images, which makes the dataset more complex and realistic. We train all of our models with the official training set and test the models on the official validation set, and we train our models from raw pixels.

### 4.2 Implementation details

**Configuration of BRN and FRN:** BRN has 8 modular routing layers, and each layer has $8/16/32$ modules. In the terminology of the ResNext paper, the cardinality is $8/16/32$ and the depth is $26 = 3 \times 8 + 2$. It is notable that we replace the 3rd batch normalization (bn) of the bottleneck block with group normalization (gn). Although most studies have shown that gn does not perform as well as bn when the batch size is large enough, and it is rarely seen that gn is applied to the ResNext model. But for BRN, the combination of gn and ResNext will greatly improve the performance of the model. We regard this as an interesting discovery, and we provide a possible explanation in the ablation study. FRN is based on ResNet34, that is, each modular layer is a *basicblock* and has 16 modular layers in total, the number of modules for each layer depends on the second convolutional layer's filter number.

**Visual Network:** all images are resized to $480 \times 360$ and we found that whether the image size is $480 \times 360$ or $224 \times 224$ has little effect on the performance of BRN, so BRN takes images of $224 \times 224$ as input. We also found that sending the question features together with the visual features into the classifier will help the BRN training more stable. In order to achieve the best performance, we concatenate $\tilde{y}_L$ with two coordinate feature maps indicating relative $x$ and $y$ spatial position.

**Textual Network:** the word embedding size is set to 200, the GRU hidden size is set to 512, we observe that the hidden size set to 512 or 1024 has little effect on the final accuracy. The parameters of the GRU and word embedding layer are initialized with orthogonal initialization and uniform initialization respectively.

**Routing Network:** although the temperature $\tau$ can be annealed to a small value during training, we find that just keeping it a constant value 1.0 can get decent accuracy. Note that we need the generation of routing path to be stochastic to exploit more possible routing paths during training, but during test phrase, we want the generation of routing path to be deterministic, so just use the sigmoid function with a threshold 0.5 to convert the $\tilde{P}$ to $P$. We also initialize the parameters of the routing network so that the probability of each module being executed at the beginning is 0.7, that is, the bias of the fc layer of routing network is initialized to $\log(\frac{0.7}{1-0.7})$.

**Training:** all the models are trained with ADAM(Kingma and Ba 2015) optimizer, betas are set to $(0.9, 0.999)$, batch-size is set to 64, learning rate is set to $3e-4$. We observe that applying a warmup scheme(Goyal et al. 2017) can help the

FRN to achieve better performance, *i.e.,* we start training our model with a small learning rate 3e-6, and slowly increase the learning rate until it reaches 3e-4, then use $3e-4$ to train the model until the end.

## 4.3 Results on CLEVR dataset

The results of all the compared methods on CLEVR are shown in Table 1, as we claimed before, the standard VQA methods like CNN+LSTM, CNN+LSTM+MCB, CNN+LSTM+SA perform poorly on this challenging dataset. Compared with the explicit reasoning methods N2NMN, PG+EE, TbD, and NS-CL, which are also based on modular networks, most of our proposed models outperform N2NMN and PG+EE, and our best-performed model is as good as NS-CL. Please note that all the four explicit methods require expert knowledge to analyze subquestion set and design corresponding modules for solving them, and all methods except NS-CL require extra supervisory information. So although it is slightly unfair to compare our models with these methods, our models still achieve competitive performance. We attribute such results to the combination of routing mechanism and modular networks. Compared with question guided methods like FiLM and QGHC, our question guided modular routing networks still outperform these methods, which validates the effectiveness of routing mechanism, especially that QGHC and BRN are both based on ResNext, but BRN outperforms QGHC by a large margin. Compared with the implicit reasoning methods like FiLM and MAC, and the prominent Relation Network, our best-performed model still outperforms most of these methods and achieves the same level of performance as MAC.

By comparing the three BRN models with different branch numbers, we can basically conclude that increasing the number of modules per layer will improve the performance accordingly. This makes sense, as the diversity of routing path can be expressed as $2^{L*M}$ (the routing path consists of $L*M$ switches, each of which has two states), if we keep the number of module layers $L$ consistent and increase the number of modules per layer $M$, the diversity of routing path will also increase. Since the diversity of routing path largely reflects the model's representation power, so increasing the number of modules per layer can enhance the performance of the model finally. We hope this discovery can provide some help for future research on routing models. By comparing the model performance with or without attention component we propose in Section 3.4 , we can find the attention component can improve the overall accuracy by about 0.7%.

Finally, on `Exist`, `Query Attribute` and `Compare Attribute` question types, our best-performed model has reached nearly 100% accuracy, higher than all other methods. But on `Compare Numbers` question type, our models do not perform well. However previous study (Perez et al. 2018) has pointed out that on `Compare Numbers` question type, model trained from feature map (extracted from a pre-trained model) will perform better than the model trained from raw pixel.

## 4.4 Results on CLEVR-Humans

To further validate the reasoning ability and generalization ability of our model, we next provide the results on CLEVR-Humans. To make a fair comparison with other methods, the reported result is fine-tuned from previous best-performed BRN and FRN models on CLEVR dataset. Also, pre-trained word embeddings are not used.

As shown in the last column of Table 1, our best-performed model outperforms all compared methods. The performance difference is more significant compared to CLEVR which validates that our model has better robustness and strong reasoning ability.

## 4.5 Ablation study of BRN

We conduct ablation studies of BRN on CLEVR dataset to analyze the influence of each decision to the model performance.

**Should we send question features to the classifier with image features?** As shown in Table 2, when the number of modules per modular layers is set to 16, sending question features together with image features to the classifier will provide 3.3% improvement over sending only image features to the classifier. However, we did not find a similar improvement in FRN, perhaps because the performance of FRN itself is high enough.

**Should we replace the 3rd normalization layer of the ResNext block from bn to gn?** First of all, why use gn instead of bn when the batchsize is already large enough? Since bn keeps $C$ pairs of mean and variance for each channel within a mini-batch. And a generally accepted rule is that the model with bn does not work well when the batchsize is too small. One possible reason is that the smaller batchsize will cause the statistics (mean and var) unstable, and change the distribution of feature map all the time. In the scenario of dynamic routing, bn's statistics will also be unstable. To validate our analysis, we replace the 3rd normalization layer of the bottleneck block from bn to gn. The results can be found in Table 2, we can see the improvements are very significant after the replacement. As for why FRN can achieve good results without the replacement, we suspect that it may be due to the higher module execution ratio of FRN (please refer to section 4.6).

**How to select the number of modules per layer?** In the previous section, we have shown that for BRN with gn, increasing the number of modules per layer will improve the performance accordingly. In Table 2, we also show that the above rule apply to BRN without gn (*i.e.,* without replacing the 3 normalization layers in the ResNext block) too.

**Will deepening or widening ResNext bring performance improvements?** The default depth and width of BRN's ResNext can be denoted as [2,2,2,2] and [64,128,256,256]. The $i$-th values of these two arrays represent the depth (*i.e.,* the number of bottleneck blocks) and width (i.e., channel number) of the $i$-th stage, respectively. The 4 stages are divided according to the downsampling operation. Based on this default configuration, we experiment with deepened (depth from [2,2,2,2] to [2,4,6,3]) or widened (width from [64,128,256,256] to [64,128,256,512]) models respectively.

From Table 2 we can find that deepening BRN will improve the accuracy by 0.3%, however widening BRN will result in a 0.3% accuracy drop. Therefore, the overall difference between the three models is not significant, but whether it is deepened or widened will increase the running time and slow down the convergence speed, so we finally use the default width and depth.

| Model | Overall | Count | Exist | Cmp. Num. | Query Attr. | Cmp. Attr. |
|---|---|---|---|---|---|---|
| **Ablation experiment on whether to feed q to the classifier** | | | | | | |
| c=16, w/o q | 89.7 | 84.7 | 95.9 | 81.9 | 90.2 | 94.6 |
| c=16, w/ q | 93.0 | 90.8 | 98.6 | 83.6 | 91.7 | 98.9 |
| **Ablation experiment on whether to replace bn with gn** | | | | | | |
| c=8, w/ q | 75.1 | 66.8 | 84.7 | 75.8 | 80.3 | 68.4 |
| c=8, w/ q, gn | 86.8 | 79.5 | 93.4 | 75.1 | 90.0 | 91.1 |
| c=16, w/ q | 93.0 | 90.8 | 98.6 | 83.6 | 91.7 | 98.9 |
| c=16, w/ q, gn | 94.7 | 90.8 | 98.1 | 82.8 | 97.3 | 98.4 |
| c=32, w/ q | 94.0 | 93.0 | 99.6 | 85.9 | 91.9 | 99.4 |
| c=32, w/ q, gn | 97.9 | 95.4 | 98.9 | 92.9 | 99.4 | 99.7 |
| **Ablation experiment on the number of modules per-layer** | | | | | | |
| c=8, w/ q | 75.1 | 66.8 | 84.7 | 75.8 | 80.3 | 68.4 |
| c=16, w/ q | 93.0 | 90.8 | 98.6 | 83.6 | 91.7 | 98.9 |
| c=32, w/ q | 94.0 | 93.0 | 99.6 | 85.9 | 91.9 | 99.4 |
| c=8, w/ q, gn | 86.8 | 79.5 | 93.4 | 75.1 | 90.0 | 91.1 |
| c=16, w/ q, gn | 94.7 | 90.8 | 98.1 | 82.8 | 97.3 | 98.4 |
| c=32, w/ q, gn | 97.9 | 95.4 | 98.9 | 92.9 | 99.4 | 99.7 |
| **Ablation experiment on increasing the depth and width** | | | | | | |
| c=32, w/ q | 94.0 | 93.0 | 99.6 | 85.9 | 91.9 | 99.4 |
| c=32, w/ q, deeper | 94.3 | 94.1 | 99.6 | 85.4 | 92.1 | 99.5 |
| c=32, w/ q, wider | 93.7 | 91.7 | 99.0 | 83.8 | 93.0 | 99.1 |

Table 2: Ablation studies of BRN on CLEVR dataset. The notation c is the abbreviation for cardinality, and represents the number of modules per modular layer; w/ q represents sending both image features and question features to the classifier; w/o q represents only sending the image features to the classifier; gn denotes using group normalization instead of batch normalization; deeper means more bottleneck blocks and wider means more channel number.

## 4.6 Visualization

**T-SNE visualization of routing paths**
To investigate what the routing network learns, we use t-SNE (Maaten and Hinton 2008) to visualize the routing paths in 2D embedding. Specifically, we divide the routing path into 4 stages according to the downsampling in the visual network, then flatten the routing path of each stage into a vector, and finally project it to the 2D space through t-SNE. For efficiency, we did not visualize the entire CLEVR dataset, but randomly select 500 instances on the validation set for each officially provided question subtype (*i.e.,* subtypes subdivided from the previous 5 types). The visualization and detailed 13 question subtypes can be seen in Figure 2. Each point represents an instance, and points of the same question subtype are labeled with the

same color. From the figure, we can discover some interesting phenomenons, we list some of them as follows: a) as the stage increases, data points with different question types can be discriminated better. b) The visualization of stage1 may be confusing at first glance, but note that data points with question subtype of query_color and equal_color are clustered together and data points with question subtype of query_material and equal_material are clustered together, this makes sense because the first few layers of CNN are responsible for detecting features about colors, textures, and edges (Olah, Mordvintsev, and Schubert 2017). c) data points with the question subtype belonging to the same question type (*e.g.,* greater_than, less_than and equal_integer belong to the Compare Numbers type) are tend to cluster together. d) data points with the question type of count and exist are clustered together, this makes sense as exist is a special case of count.

**How many modules are executed for each layer**
To investigate how many modules are executed in each layer, we provide the box plot (Figure 3) of the module execution ratio (*i.e.,* the number of executed modules divided by the total number of the modules) for each layer of FRN and BRN(c=32) on the validation set. First of all, it can be observed that both FRN (refer to Figure 3a) and BRN (refer to Figure 3b) follow a rule: the execution ratio becomes lower as layers going deeper, this makes sense as the features extracted at early layers are generic, and the features extracted at higher layers are more related to the VQA task. Next we can observe that the overall execution ratio of the FRN model is not low, however please note that with the reduction of filters, current deep learning frameworks cannot reduce the running time accordingly, so we did not pay much attention to it. Finally, note that even if we do not add any sparsity constraints, BRN still has a considerable reduction of about **43%**. We believe that this number can be increased with the addition of sparsity constraints. And with the further support of the deep learning framework, the routing model will have more useful applications in the future.

**Saliency map**
To investigate which pixels affect the prediction most, we visualize the saliency map of the same image to different questions. Concretely, we calculate the gradient of each pixel with respect to the max prediction score. Figure 4 quantitatively shows that QGMRN can identify the key pixels for prediction. More examples are available in appendix A.

## 5 Conclusion

In this paper, we propose the Question Guided Modular Routing Networks for VQA which can fuse the visual and textual modalities at multiple semantic levels and learn to reason by routing between the generic modules, different interesting variant models can be generated by changing the granularity of the module. In the experiments, we show that our models have achieved the state-of-the-art performance. From the perspective of routing models, we find a suitable application for them where static models struggle, we also successfully applied them to large models and large-scale datasets. In particular, we provide two interesting findings:
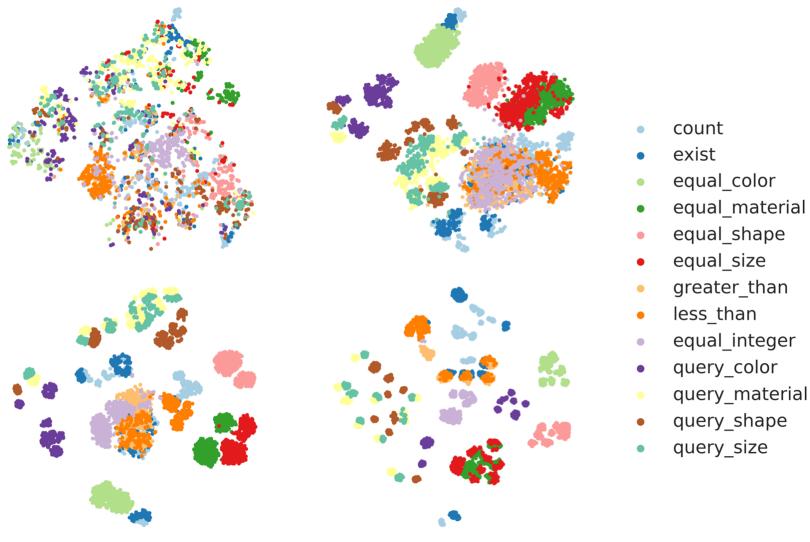
(a) FRN



(b) BRN(c=32)

Figure 2: The t-SNE visualization of routing paths of stage1, stage2, stage3 and stage4 are represented at top-left, top-right, bottom-left and bottom-right sub-figures respectively. Low-level question subtypes are clustered in the earlier stages and high-level question subtypes are clustered in the later stages. Best seen on the computer, in colour and zoomed in.

Figure 3: The module execution ratio in each layer. The overall execution ratio of FRN is above BRN.



Q: What color is the rubber cube?
A: red

Q: How many brown metal cylinders have the same size as the purple cylinder?
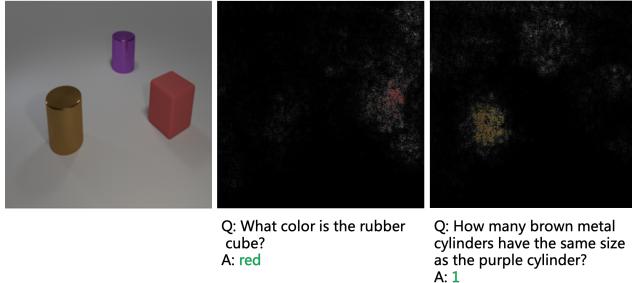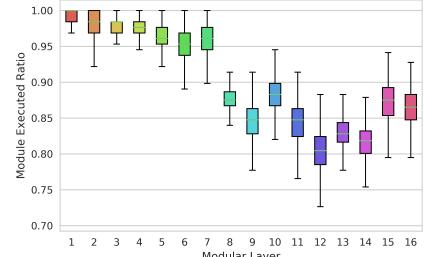A: 1

Figure 4: The saliency map indicates the key pixels for reasoning with respect to different questions.
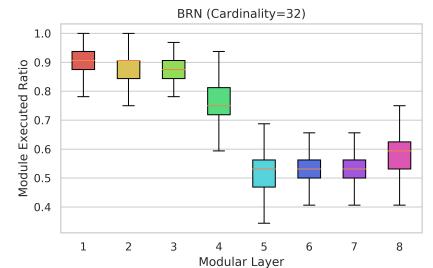
1. The performance of the model is positively correlated with the number of modules per layer; 2. When the execution ratio is low, using gn will be better than bn. We believe that routing models will play an important role in future multimodal fusion and embedding.

## References

[Anderson et al. 2018] Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; and Zhang, L. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*.

[Andreas et al. 2016] Andreas, J.; Rohrbach, M.; Darrell, T.; and Klein, D. 2016. Neural module networks. In *CVPR*.

[Antol et al. 2015] Antol, S.; Agrawal, A.; Lu, J.; Mitchell, M.; Batra, D.; Lawrence Zitnick, C.; and Parikh, D. 2015. Vqa: Visual question answering. In *ICCV*.

[Ben-Younes et al. 2017] Ben-Younes, H.; Cadene, R.; Cord, M.; and Thome, N. 2017. Mutan: Multimodal tucker fusion for visual question answering. In *ICCV*.

[Bengio et al. 2016] Bengio, E.; Bacon, P.-L.; Pineau, J.; and Precup, D. 2016. Conditional computation in neural networks for faster models. In *ICLR workshop*, volume abs/1511.06297.

[Bengio, Léonard, and Courville 2013] Bengio, Y.; Léonard, N.; and Courville, A. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR* abs/1308.3432.

[Cho et al. 2014] Cho, K.; Van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; and Bengio, Y. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*.

[De Vries et al. 2017] De Vries, H.; Strub, F.; Mary, J.; Larochelle, H.; Pietquin, O.; and Courville, A. C. 2017. Modulating early visual processing by language. In *NIPS*, 6594–6604.

[Eigen, Ranzato, and Sutskever 2013] Eigen, D.; Ranzato, M.; and Sutskever, I. 2013. Learning factored representations in a deep mixture of experts. In *ICLR*.

[Fukui et al. 2016] Fukui, A.; Park, D. H.; Yang, D.; Rohrbach, A.; Darrell, T.; and Rohrbach, M. 2016. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*.

[Gao et al. 2018] Gao, P.; Li, H.; Li, S.; Lu, P.; Li, Y.; Hoi, S. C.; and Wang, X. 2018. Question-guided hybrid convolution for visual question answering. In *ECCV*.

[Goyal et al. 2017] Goyal, P.; Dollár, P.; Girshick, R.; Noordhuis, P.; Wesolowski, L.; Kyrola, A.; Tulloch, A.; Jia, Y.; and He, K. 2017. Accurate, large minibatch sgd: Training imagenet in 1 hour. *CoRR* abs/1706.02677.

[He et al. 2016] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*, 770–778.

[Hochreiter and Schmidhuber 1997] Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.

[Hu et al. 2017] Hu, R.; Andreas, J.; Rohrbach, M.; Darrell, T.; and Saenko, K. 2017. Learning to reason: End-to-end module networks for visual question answering. In *ICCV*.

[Hudson and Manning 2018] Hudson, D. A., and Manning, C. D. 2018. Compositional attention networks for machine reasoning. In *ICLR*.

[Ioffe and Szegedy 2015] Ioffe, S., and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR* abs/1502.03167.

[Jacobs et al. 1991] Jacobs, R. A.; Jordan, M. I.; Nowlan, S. J.; Hinton, G. E.; et al. 1991. Adaptive mixtures of local experts. *Neural computation* 3(1):79–87.

[Jang, Gu, and Poole 2017] Jang, E.; Gu, S.; and Poole, B. 2017. Categorical reparameterization with gumbel-softmax. In *ICLR*.

[Johnson et al. 2017a] Johnson, J.; Hariharan, B.; van der Maaten, L.; Fei-Fei, L.; Lawrence Zitnick, C.; and Girshick, R. 2017a. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2901–2910.

[Johnson et al. 2017b] Johnson, J.; Hariharan, B.; van der Maaten, L.; Hoffman, J.; Fei-Fei, L.; Lawrence Zitnick, C.; and Girshick, R. 2017b. Inferring and executing programs for visual reasoning. In *ICCV*.

[Jordan and Jacobs 1994] Jordan, M. I., and Jacobs, R. A. 1994. Hierarchical mixtures of experts and the em algorithm. *Neural computation* 6(2):181–214.

[Kazemi and Elqursh 2017] Kazemi, V., and Elqursh, A. 2017. Show, ask, attend, and answer: A strong baseline for visual question answering. *CoRR* abs/1704.03162.

[Kim et al. 2016] Kim, J.-H.; Lee, S.-W.; Kwak, D.-H.; Heo, M.-O.; Kim, J.; Ha, J.-W.; and Zhang, B.-T. 2016. Multimodal residual learning for visual qa. In *NIPS*.

[Kingma and Ba 2015] Kingma, D. P., and Ba, J. 2015. Adam: A method for stochastic optimization. In *ICLR*.

[Kirsch, Kunze, and Barber 2018] Kirsch, L.; Kunze, J.; and Barber, D. 2018. Modular networks: Learning to decompose neural computation. In *NeurIPS*.

[Maaten and Hinton 2008] Maaten, L. v. d., and Hinton, G. 2008. Visualizing data using t-sne. *Journal of machine learning research* 9(Nov):2579–2605.

[Maddison, Mnih, and Teh 2017] Maddison, C. J.; Mnih, A.; and Teh, Y. W. 2017. The concrete distribution: A continuous relaxation of discrete random variables. In *ICLR*.

[Malinowski, Rohrbach, and Fritz 2015] Malinowski, M.; Rohrbach, M.; and Fritz, M. 2015. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 1–9.

[Mao et al. 2019] Mao, J.; Gan, C.; Kohli, P.; Tenenbaum, J. B.; and Wu, J. 2019. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *ICLR*.

[Mascharka et al. 2018] Mascharka, D.; Tran, P.; Soklaski, R.; and Majumdar, A. 2018. Transparency by design: Closing the gap between performance and interpretability in visual reasoning. In *CVPR*.

[Olah, Mordvintsev, and Schubert 2017] Olah, C.; Mordvintsev, A.; and Schubert, L. 2017. Feature visualization. *Distill*. https://distill.pub/2017/feature-visualization.

[Perez et al. 2018] Perez, E.; Strub, F.; De Vries, H.; Dumoulin, V.; and Courville, A. 2018. Film: Visual reasoning with a general conditioning layer. In *AAAI*.

[Ramachandran and Le 2019] Ramachandran, P., and Le, Q. V. 2019. Diversity and depth in per-example routing models. In *ICLR*.

[Rosenbaum, Klinger, and Riemer 2018] Rosenbaum, C.; Klinger, T.; and Riemer, M. 2018. Routing networks: Adaptive selection of non-linear functions for multi-task learning. In *ICLR*.

[Santoro et al. 2017] Santoro, A.; Raposo, D.; Barrett, D. G.; Malinowski, M.; Pascanu, R.; Battaglia, P.; and Lillicrap, T. 2017. A simple neural network module for relational reasoning. In *NIPS*.

[Shazeer et al. 2017] Shazeer, N.; Mirhoseini, A.; Maziarz, K.; Davis, A.; Le, Q. V.; Hinton, G. E.; and Dean, J. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *ICLR*.

[Simonyan and Zisserman 2014] Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556.

[Vaswani et al. 2017] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *NIPS*.

[Veit and Belongie 2018] Veit, A., and Belongie, S. J. 2018. Convolutional networks with adaptive inference graphs. In *ECCV*.

[Wu et al. 2018a] Wu, C.; Liu, J.; Wang, X.; and Dong, X. 2018a. Object-difference attention: A simple relational attention for visual question answering. In *2018 ACM Multimedia Conference on Multimedia Conference*, 519–527. ACM.

[Wu et al. 2018b] Wu, Z.; Nagarajan, T.; Kumar, A.; Rennie, S.; Davis, L. S.; Grauman, K.; and Feris, R. S. 2018b. Blockdrop: Dynamic inference paths in residual networks. *CVPR*.

[Xie et al. 2017] Xie, S.; Girshick, R. B.; Dollár, P.; Tu, Z.; and He, K. 2017. Aggregated residual transformations for deep neural networks. *CVPR* 5987–5995.

[Yang et al. 2016] Yang, Z.; He, X.; Gao, J.; Deng, L.; and Smola, A. 2016. Stacked attention networks for image question answering. In *CVPR*, 21–29.

[Yellott Jr 1977] Yellott Jr, J. I. 1977. The relationship between luce's choice axiom, thurstone's theory of comparative judgment, and the double exponential distribution. *Journal of Mathematical Psychology* 15(2):109–144.

[Zhou et al. 2015] Zhou, B.; Tian, Y.; Sukhbaatar, S.; Szlam, A.; and Fergus, R. 2015. Simple baseline for visual question answering. *CoRR* abs/1512.02167.

# A    Saliency map examples

We provide several saliency maps to illustrate the key pixels conditioned on the question in Figure 5.

# B    Pseudocode of routing algorithm

---

**Algorithm 1** Pseudocode of routing algorithm

---

**Require:** question feature $\mathbf{q}$; shape of routing path $L$ and $M$; temperature $\tau$

**Ensure:** routing path $\mathbf{P}$

    /* Sample from Logistic distribution */

1: First Sample a random tensor from Uniform distribution $\mathbf{u} \sim \text{Uniform}([0,1]^{L \times M})$, where $\mathbf{u} \in \mathbb{R}^{L \times M}$

2: Then compute $\mathbf{t} \leftarrow \log(\mathbf{u}) - \log(1 - \mathbf{u})$, and $\mathbf{t} \sim \text{Logistic}([0,1]^{L \times M})$

    /* Compute routing path */

3: First use a fc layer to get the real-valued routing path $\tilde{\mathbf{P}} \leftarrow \text{FC}(\mathbf{q})$

4: Then use Gumbel-Sigmoid trick to compute random tensor $\tilde{\mathbf{B}} \leftarrow \text{Sigmoid}(\frac{\mathbf{t} + \tilde{\mathbf{P}}}{\tau})$

5: Threshold $\tilde{\mathbf{B}}$ to get the $\tilde{\mathbf{B}}_{\text{binary}} \leftarrow \tilde{\mathbf{B}} > 0.5$

6: Finally use Straight-Through (ST) trick (Bengio, Léonard, and Courville 2013) to compute discrete routing path $\mathbf{P} \leftarrow \tilde{\mathbf{B}}_{\text{binary}} - \text{detach}(\tilde{\mathbf{B}}) + \tilde{\mathbf{B}}$

---

Q: What color is the rubber cube?
A: red

Q: How many brown metal cylinders have the same size as the purple cylinder?
A: 1

Q: The tiny shiny cylinder has what color?
A: Brown

Q: Is the color of the metal block that is right of the yellow rubber object the same as the large metal cylinder?
A: Yes

Q: What is the large red cylinder made of?
A: Metal

Q: What shape is the big metal thing that is the same color as the small cylinder?
A: Cylinder

Q: There is a thing that is in front of the brown thing; what is its color?
A: Cyan

Q: There is a large cyan matte object; does it have the same shape as the shiny object to the left of the big cylinder?
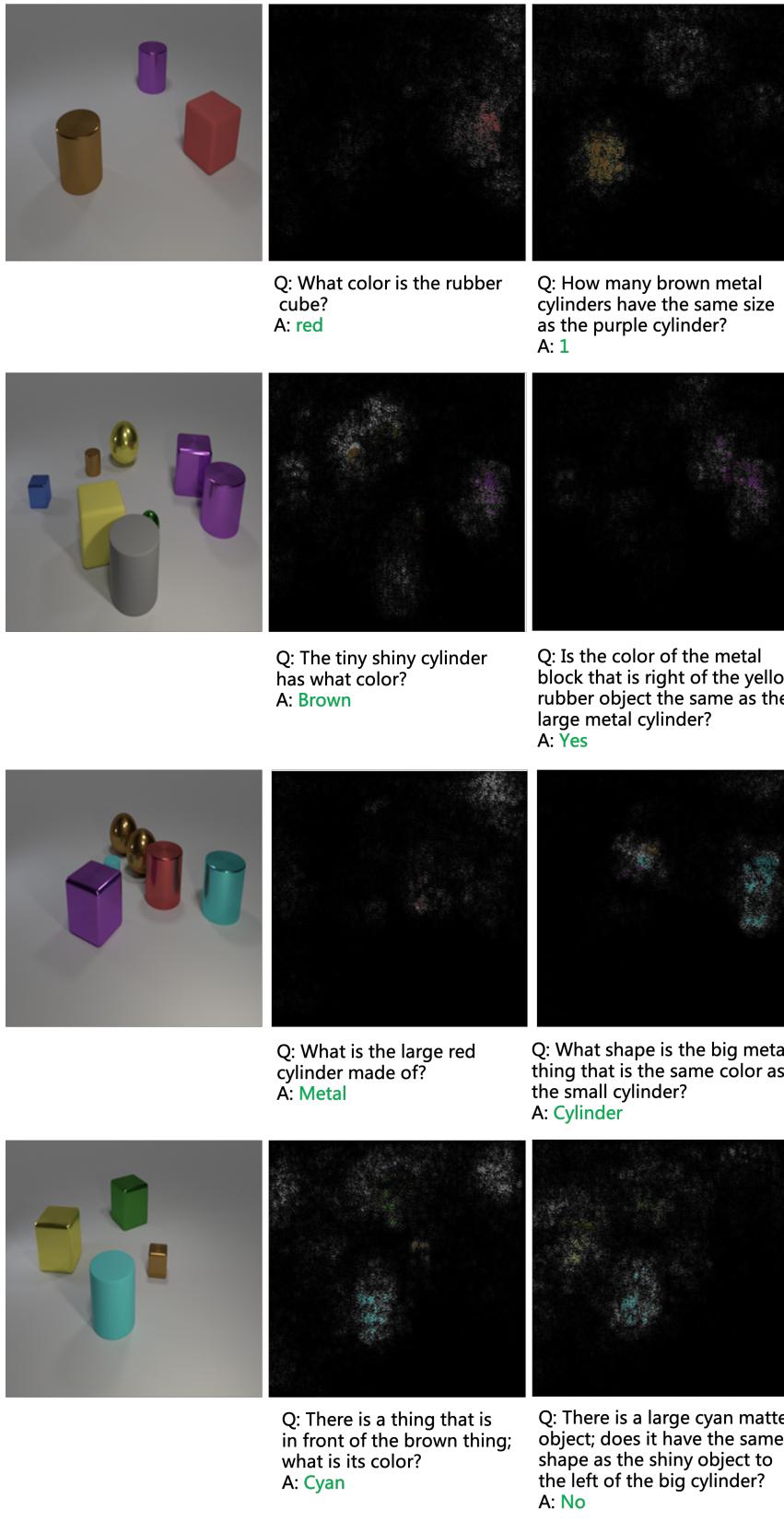A: No

Figure 5: More saliency map examples with respect to different questions.