# Explainable and Explicit Visual Reasoning over Scene Graphs

Jiaxin Shi[1*]      Hanwang Zhang[2]      Juanzi Li[1]

[1]Tsinghua University      [2]Nanyang Technological University

shijx12@gmail.com; hanwangzhang@ntu.edu.sg; lijuanzi@tsinghua.edu.cn

## Abstract

*We aim to dismantle the prevalent black-box neural architectures used in complex visual reasoning tasks, into the proposed eXplainable and eXplicit Neural Modules (XNMs), which advance beyond existing neural module networks towards using scene graphs — objects as nodes and the pairwise relationships as edges — for explainable and explicit reasoning with structured knowledge. XNMs allow us to pay more attention to teach machines how to "think", regardless of what they "look". As we will show in the paper, by using scene graphs as an inductive bias, 1) we can design XNMs in a concise and flexible fashion, i.e., XNMs merely consist of 4 meta-types, which significantly reduce the number of parameters by 10 to 100 times, and 2) we can explicitly trace the reasoning-flow in terms of graph attentions. XNMs are so generic that they support a wide range of scene graph implementations with various qualities. For example, when the graphs are detected perfectly, XNMs achieve 100% accuracy on both CLEVR and CLEVR CoGenT, establishing an empirical performance upper-bound for visual reasoning; when the graphs are noisily detected from real-world images, XNMs are still robust to achieve a competitive 67.5% accuracy on VQAv2.0, surpassing the popular bag-of-objects attention models without graph structures.*

## 1. Introduction

The prosperity of A.I. — mastering super-human skills in game playing [23], speech recognition [1], and image recognition [8, 21] — is mainly attributed to the "winning streak" of *connectionism*, more specifically, the deep neural networks [16], over the "old-school" *symbolism*, where their controversy can be dated back to the birth of A.I. in 1950s [19]. With massive training data and powerful computing resources, the key advantage of deep neural networks is the end-to-end design that generalizes to a large spec-
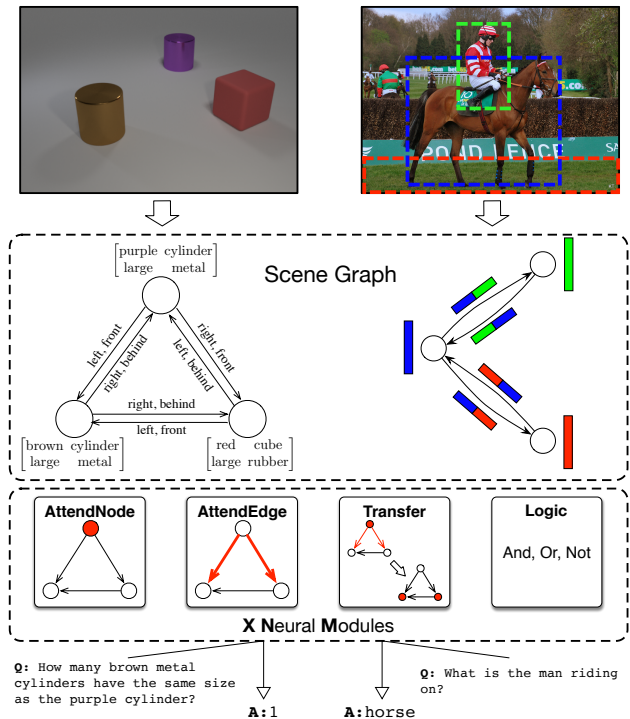


Figure 1: The flowchart of using the proposed XNMs reasoning over scene graphs, which can be represented by detected one-hot class labels (left) or RoI feature vectors (colored bars on the right). Feature colors are consistent with the bounding box colors. XNMs have 4 meta-types. Red nodes or edges indicate attentive results. The final module assembly can be obtained by training an off-the-shelf sequence-to-sequence program generator [13].

trum of domains, minimizing the human efforts in domain-specific knowledge engineering. However, large gaps between human and machines can be still observed in "high-level" vision-language tasks such as visual Q&A [4, 6, 12], which inherently requires composite reasoning (cf. Figure 1). In particular, recent studies show that the end-to-end models are easily optimized to learn the dataset "shortcut bias" but not reasoning [12].

---

Neural module networks (NMNs) [3, 12, 10, 18, 9, 27] show a promising direction in conferring reasoning ability for the end-to-end design by learning to compose the networks on-demand from the language counterpart, which implies the logical compositions. Take the question *"How many objects are left of the red cube?"* as an example, we can program the reasoning path into a composition of functional modules [18]: Attend[cube], Attend[red], Relate[left], and Count, and then execute them with the input image. We attribute the success of NMNs to the eXplainable and eXplicit (dubbed **X**) language understanding. By explicitly parsing the question into an explainable module assembly, NMNs effectively prevent the language-to-reasoning shortcut, which are frequent when using the implicit fused question representations [4, 6] (*e.g.*, the answer can be directly inferred according to certain language patterns).

However, the vision-to-reasoning shortcut still exists as an obstacle on the way of NMNs towards the real X visual reasoning. This is mainly because that the visual perception counterpart is still attached to reasoning [18], which is inevitably biased to certain vision patterns. For example, on the CLEVR CoGenT task, which provides novel object attributes to test the model's generalization ability (*e.g.*, cubes are blue in the training set but red in the test set), we observe significant performance drop of existing NMNs [13, 18] (*e.g.*, red cubes in the test set cannot be recognized as "cube"). Besides, the reusability of the current module design is limited. For example, the network structure of the Relate module in [18] must be carefully designed using a series of dilated convolutions to achieve good performance. Therefore, how to design a complete inventory of X modules is still an tricky engineering.

In this paper, we advance NMN towards X visual reasoning by using the proposed eXplainable and eXplicit **N**eural **M**odules (XNMs) reasoning over *scene graphs*. By doing this, we can insulate the "low-level" visual perception from the modules, and thus can prevent reasoning shortcut of both language and vision counterpart. As illustrated in Figure 1, a scene graph is the knowledge representation of a visual input, where the nodes are the entities (*e.g.*, cylinder, horse) and the edges are the relationships between entities (*e.g.*, left, ride). In particular, we note that scene graph detection *per se* is still a challenging task in computer vision [29], therefore, we allow XNMs to accept scene graphs with different detection qualities. For example, the left-hand side of Figure 1 is one extreme when the visual scene is clean and closed-vocabulary, *e.g.*, in CLEVR [12], we can have almost perfect scene graphs where the nodes and edges can be represented by one-hot class labels; the right-hand side shows another extreme when the scene is cluttered and open-vocabulary in practice, the best we have might be merely a set of object proposals. Then, the nodes are RoI features and the edges are their concatenations.

Thanks to scene graphs, our XNMs only have 4 meta-types: 1) AttendNode, finding the queried entities, 2) AttendEdge, finding the queried relationships, 3) Transfer, transforming the node attentions along the attentive edges, and 4) Logic, performing basic logical operations on attention maps. All types are fully X as their outputs are pure graph attentions that are easily traceable and visible. Moreover, these meta modules are only specific to the generic graph structures, and are highly reusable to constitute different composite modules for more complex functions. For example, we do not need to carefully design the internal implementation details for the module Relate as in [18]; instead, we only need to combine AttendEdge and Transfer in XNMs.

We conduct extensive experiments [1] on two visual Q&A benchmarks and demonstrate the following advantages of using XNMs reasoning over scene graphs:

1. We achieve 100% accuracy by using the ground-truth scene graphs and programs on both CLEVR [12] and CLEVR-CoGent, revealing the performance upper-bound of XNMs, and the benefits of disentangling "high-level" reasoning from "low-level" perception.
2. Our network requires significantly less parameters while achieves better performance than previous state-of-the-art neural module networks, due to the conciseness and high-reusability of XNMs.
3. XNMs are flexible to different graph qualities, *e.g.*, it achieves competitive accuracy on VQAv2.0 [6] when scene graphs are noisily detected.
4. We show qualitative results to demonstrate that our XNMs reasoning is highly explainable and explicit.

## 2. Related Work

**Visual Reasoning.** It is the process of analyzing visual information and solving problems based on it. The most representative benchmark of visual reasoning is CLEVR [12], a diagnostic visual Q&A dataset for compositional language and elementary visual reasoning. The majority of existing methods on CLEVR can be categorized into two families: 1) holistic approaches [12, 22, 20, 11], which embed both the image and question into a feature space and infer the answer by feature fusion; 2) neural module approaches [3, 10, 13, 18, 9, 27], which first parse the question into a program assembly of neural modules, and then execute the modules over the image features for visual reasoning. Our XNM belongs to the second one but replaces the visual feature input with scene graphs.

**Neural Module Networks.** They dismantle a complex question into several sub-tasks, which are easier to answer and more transparent to follow the intermediate outputs.

---

[1]Our codes are public at https://github.com/shijx12/XNM-Net

Modules are pre-defined neural networks that implement the corresponding functions of sub-tasks, and then are assembled into a layout dynamically, usually by a sequence-to-sequence program generator given the input question. The assembled program is finally executed for answer prediction [10, 13, 18]. In particular, the program generator is trained based on the human annotations of desired layout or with the help of reinforcement learning due to the non-differentiability of layout selection. Recently, Hu *et al*. [9] proposed StackNMN, which replaces the hard-layout with soft and continuous module layout and performs well even without layout annotations at all. Our XNM experiments on VQAv2.0 follows their soft-program generator.

Recently, NS-VQA [27] firstly built the reasoning over the object-level structural scene representation, improving the accuracy on CLEVR from the previous state-of-the-art 99.1% [18] to an almost perfect 99.8%. Their scene structure consists of objects with detected labels, but lacked the relationships between objects, which limited its application on real-world datasets such as VQAv2.0 [6]. In this paper, we propose a much more generic framework for visual reasoning over scene graphs, including object nodes and relationship edges represented by either labels or visual features. Our scene graph is more flexible and more powerful than the table structure of NS-VQA.

**Scene Graphs.** This task is to produce graph representations of images in terms of objects and their relationships. Scene graphs have been shown effective in boosting several vision-language tasks [14, 25, 28, 5]. To the best of our knowledge, we are the first to design neural module networks that can reason over scene graphs. However, scene graph detection is far from satisfactory compared to object detection [26, 29, 17]. To this end, our scene graph implementation also supports cluttered and open-vocabulary in real-world scene graph detection, where the nodes are merely RoI features and the edges are their concatenations.

# 3. Approach

We build our neural module network over scene graphs to tackle the visual reasoning challenge. As shown in Figure 2, given an input image and a question, we first parse the image into a scene graph and parse the question into a module program, and then execute the program over the scene graph. In this paper, we propose a set of generic base modules that can conduct reasoning over scene graphs — eXplainable and eXplicit Neural Modules (XNMs) — as the reasoning building blocks. We can easily assemble these XNMs to form more complex modules under specific scenarios. Besides, our XNMs are totally attention-based, making all the intermediate reasoning steps transparent.

## 3.1. Scene Graph Representations

We formulate the scene graph of an image as $(\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{\mathbf{v}_1, \cdots, \mathbf{v}_N\}$ are graph nodes corresponding to $N$ detected objects, and $\mathbf{v}_i$ denotes the feature representation of the $i$-th object. $\mathcal{E} = \{\mathbf{e}_{ij}|i, j = 1, \cdots, N\}$ are graph edges corresponding to relations between each object pairs, and $\mathbf{e}_{ij}$ denotes the feature representation of the relation from object $i$ to object $j$ (Note that edges are directed).

Our XNMs are generic for scene graphs of different quality levels of detection. We consider two extreme settings in this paper. The first is the ground-truth scene graph with labels, denoted by **GT**, that is, using ground-truth objects as nodes, ground-truth object label embeddings as node features, and ground-truth relation label embeddings as edge features. In this setting, scene graphs are annotated with fixed-vocabulary object labels and relationship labels, *e.g.*, defined in CLEVR dataset [12]. We collect all the $C$ labels into a dictionary, and use an embedding matrix $\mathbf{D} \in \mathbb{R}^{C \times d}$ to map a label into a $d$-dimensional vector. We represent the nodes and edges using the concatenation of their corresponding label embeddings.

The second setting is totally detected and label-agnostic, denoted by **Det**, that is, using detected objects as nodes, RoI visual features as node features, and the fusion of two node features as the edge features. For example, the edge features can be represented by concatenating the two related node features, *i.e.*, $\mathbf{e}_{ij} = [\mathbf{v}_i; \mathbf{v}_j]$. As an another example, in CLEVR where the edges are only about spatial relationships, we use the difference between detected coordinates of object pairs as the edge embedding. More details are in Section 4.

We use the GT setting to demonstrate the performance upper-bound of our approach when a perfect scene graph detector is available along with the rapid development of visual recognition, and use the Det setting to demonstrate the practicality in open domains.

## 3.2. X Neural Modules

As shown in Figure 1, our XNMs have four meta-types and are totally attention-based. We denote the node attention weight vector by $\mathbf{a} \in [0, 1]^N$ and the weight of the $i$-th node by $a_i$. The edge attention weight matrix is denoted by $\mathbf{W} \in [0, 1]^{N \times N}$, where $W_{ij}$ represents the weight of edge from node $i$ to node $j$.

**AttendNode[query].** This most basic and intuitive operation is to find the relevant objects given an input query (*e.g.*, find all ["cubes"]). For the purpose of semantic computation, we first encode the query into a vector $\mathbf{q}$. This X module takes the query vector as input, and produces the node attention vector by the following function:

$$\mathbf{a} = f(\mathcal{V}, \mathbf{q}). \tag{1}$$

The implementation of $f$ is designed according to a specific scene graph representation, as long as $f$ is differentiable and $\text{range}(f) = [0, 1]$.
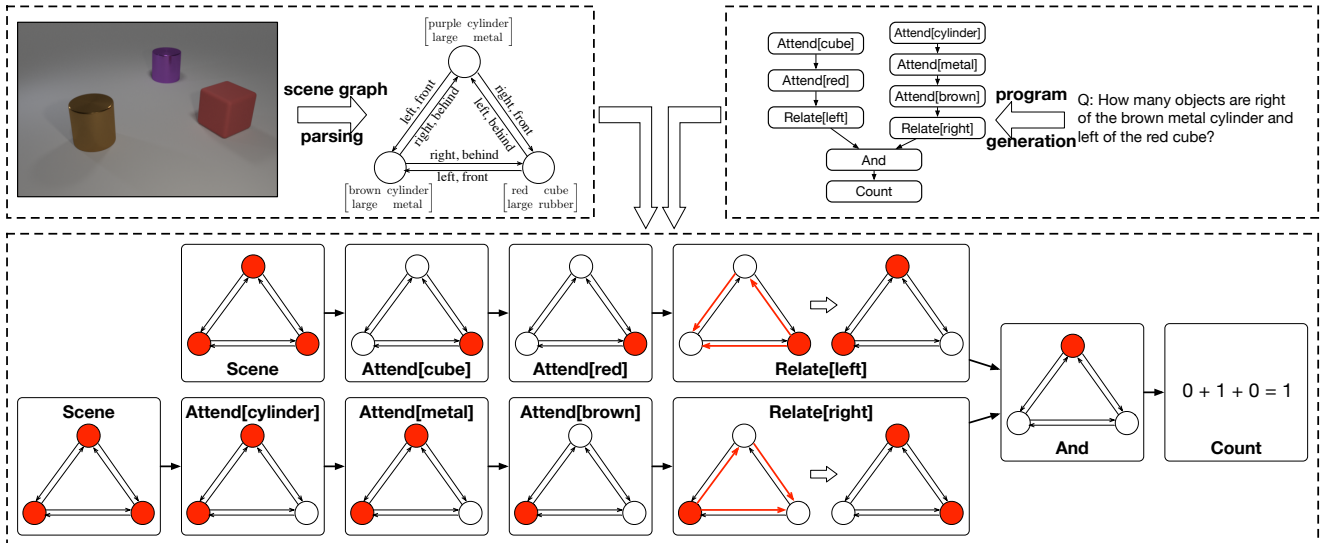
Figure 2: To answer a question about an image, we need to 1) parse the image into a scene graph, 2) parse the question into a module program, and 3) reasoning over the scene graph. Here, we show the reasoning details of an example from CLEVR. The nodes and edges in red are attended. `Scene` is a dummy placeholder module that attends all nodes. All intermediate steps of our XNMs are explainable and explicit.

**AttendEdge[query].** Though object attention is a widely-used mechanism for better visual understanding, it is unable to capture the interaction between objects and thus is weak in the complex visual reasoning [30]. This X module aims to find the relevant edges given an input query (*e.g.*, find all edges that are ["left"]). After encoding the query into $\mathbf{q}$, we compute the edge attention matrix by the following function:

$$\mathbf{W} = g(\mathcal{E}, \mathbf{q}), \qquad (2)$$

where $g$ is defined according to a specific scene graph representation, as long as $g$ is differentiable and $\text{range}(g) = [0, 1]$.

**Transfer.** With the node attention vector $\mathbf{a}$ and the edge attention matrix $\mathbf{W}$, we can transfer the node weights along the attentive relations to find new objects (*e.g.*, find objects that are ["left"] to the ["cube"]). Thanks to the graph structure, to obtain the updated node attention $\mathbf{a}'$, we merely need to perform a simple matrix multiplication:

$$\mathbf{a}' = \text{norm}(\mathbf{W}^\top \mathbf{a}), \qquad (3)$$

where norm assert the values in $[0, 1]$ by dividing the maximum value if any entry exceeds 1. Here, $W_{ij}$ indicates how many weights will flow from object $i$ to object $j$, and $a'_i = \sum_{j=1}^{N} W_{ji} a_j$ is the total received weights of object $i$. This module reallocates node attention in an efficient and fully-differentiable manner.

**Logic.** Logical operations are crucial in complex reasoning cases. In XNM, logical operations are performed on one or more attention weights to produce a new attention. We

define three logical X modules: **And**, **Or**, and **Not**. Without loss of generality, we discuss all these logical modules on node attention vectors, and the extension to edge attention is similar. The And and Or modules are binary, that is, take two attentions as inputs, while the Not module is unary. The implementation of these logical X modules are as follows:

$$\begin{aligned} \text{And}(\mathbf{a}^1, \mathbf{a}^2) &= \min(\mathbf{a}^1, \mathbf{a}^2), \text{Not}(\mathbf{a}) = 1 - \mathbf{a}, \\ \text{Or}(\mathbf{a}^1, \mathbf{a}^2) &= \max(\mathbf{a}^1, \mathbf{a}^2). \end{aligned} \qquad (4)$$

These four meta-types of XNMs constitute the base of our graph reasoning. They are explicitly executed on attention maps, and all intermediate results are explainable. Besides, these X modules are totally differentiable. We can flexibly assemble them into composite modules for more complex functions, which can be still trained end-to-end.

### 3.3. Implementations

To apply XNMs in practice, we need to consider these questions: **(1)** How to implement the attention functions $f, g$ in Eq. (1) and Eq. (2)? **(2)** How to compose our X modules into composite reasoning modules? **(3)** How to predict the answer according to the attentive results? **(4)** How to parse the input question to an executable module program?

#### 3.3.1 Attention Functions

We use different attention functions for different scene graph settings. In the GT setting, as annotated labels are mostly mutually exclusive (*e.g.*, "red" and "green"), we compute the node attention using the softmax function over

4

the label space. Specifically, given a query vector $\mathbf{q} \in \mathbb{R}^d$, we first compute its attention distribution over all labels by $\mathbf{b} = \text{softmax}(\mathbf{D} \cdot \mathbf{q})$, where $\text{length}(\mathbf{b}) = C$ and $b_c$ represents the weight of the $c$-th label. Then we capture the node and edge attention by summing up corresponding label weights:

$$a_i = f(\mathcal{V}, \mathbf{q})_i = \sum_{c \in \mathcal{C}_i} b_c, \;\; W_{ij} = g(\mathcal{E}, \mathbf{q})_{ij} = \sum_{c \in \mathcal{C}_{ij}} b_c, \tag{5}$$

where $\mathcal{C}_i$ and $\mathcal{C}_{ij}$ denote the (multi-) labels of node $i$ and edge $ij$ respectively.

In the Det setting, we use the sigmoid function to compute the attention weights. Given the query $\mathbf{q} \in \mathbb{R}^d$, the node and edge attentions are:

$$a_i = f(\mathcal{V}, \mathbf{q})_i = \text{sigmoid}\left(\text{MLP}(\mathbf{v}_i)^\top \mathbf{q}\right),$$
$$W_{ij} = g(\mathcal{E}, \mathbf{q})_{ij} = \text{sigmoid}\left(\text{MLP}(\mathbf{e}_{ij})^\top \mathbf{q}\right), \tag{6}$$

where the MLP maps $\mathbf{v}_i$ and $\mathbf{e}_{ij}$ to the dimension $d$.

### 3.3.2 Composite Reasoning Modules

We list our composite reasoning modules and their implementations (*i.e.*, how they are composed by basic X modules) in the top section of Table 1. For example, `Same` module is to find other objects that have the same attribute value as the input objects (*e.g.*, find other objects with the same ["color"]). In particular, `Describe` used in `Same` is to obtain the corresponding attribute value (*e.g.*, describe one object's ["color"]), and will be introduced in the following section.

### 3.3.3 Feature Output Modules

Besides the above reasoning modules, we also need another kind of modules to map the intermediate attention to a hidden embedding $\mathbf{h}$ for feature representation, which is fed into a softmax layer to predict the final answer, or into some modules for further reasoning. We list our output modules in the bottom section of Table 1. `Exist` and `Count` sum up the node attention weights to answer yes/no and counting questions. `Compare` is for attribute or number comparisons, which takes two hidden features as inputs. **`Describe[query]`** is to transform the attentive node features to an embedding that describes the specified attribute value (*e.g.*, what is the ["color"] of attended objects).

To implement the `Describe` module, we first obtain the "raw" attentive node feature by

$$\bar{\mathbf{v}} = \sum_{i=1}^{N} a_i \mathbf{v}_i \Big/ \sum_{i=1}^{N} a_i, \tag{7}$$

and then project it into several "fine-grained" sub-spaces — describing different attribute aspects such as color and

shape — using different transformation matrices. Specifically, we define $K$ projection matrices $\mathbf{M}_1, \cdots, \mathbf{M}_K$ to map $\bar{\mathbf{v}}$ into different aspects (*e.g.*, $\mathbf{M}_1\bar{\mathbf{v}}$ represents the color, $\mathbf{M}_2\bar{\mathbf{v}}$ represents the shape, *etc.*), where $K$ is a hyper-parameter related to the specific scene graph vocabulary. The output feature is computed by

$$\texttt{Describe}(\mathbf{a}, \mathbf{q}) = \sum_{k=1}^{K} c_k(\mathbf{M}_k \bar{\mathbf{v}}), \tag{8}$$

where $\mathbf{c} = \text{Softmax}(\text{MLP}(\mathbf{q}))$ represents a probability distribution over these $K$ aspects, and $c_k$ denotes the $k$-th probability. The mapping matrixes can be learned end-to-end automatically.

Table 1: Our composite modules (the top section) and output modules (the bottom section). $\text{MLP}()$ consists of several linear and ReLU layers.

| Modules | In $\to$ Out | Implementation |
|---|---|---|
| Intersect | $\mathbf{a}^1, \mathbf{a}^2 \to \mathbf{a}'$ | And($\mathbf{a}^1, \mathbf{a}^2$) |
| Union | $\mathbf{a}^1, \mathbf{a}^2 \to \mathbf{a}'$ | Or($\mathbf{a}^1, \mathbf{a}^2$) |
| Filter | $\mathbf{a}, \mathbf{q} \to \mathbf{a}'$ | And($\mathbf{a}$, AttendNode($\mathbf{q}$)) |
| Same | $\mathbf{a}, \mathbf{q} \to \mathbf{a}'$ | Filter(Not($\mathbf{a}$), Describe($\mathbf{a}, \mathbf{q}$)) |
| Relate | $\mathbf{a}, \mathbf{q} \to \mathbf{a}'$ | Transfer($\mathbf{a}$, AttendEdge($\mathbf{q}$)) |
| Exist | $\mathbf{a} \to \mathbf{h}$ | $\text{MLP}(\sum_i a_i)$ |
| Count | | |
| Compare | $\mathbf{h}^1, \mathbf{h}^2 \to \mathbf{h}'$ | $\text{MLP}(\mathbf{h}^1 - \mathbf{h}^2)$ |
| Describe | $\mathbf{a}, \mathbf{q} \to \mathbf{h}$ | Eq. (8) |

### 3.3.4 Program Generation & Training

For datasets that have ground-truth program annotations (*e.g.*, CLEVR), we directly learn an LSTM sequence-to-sequence model [24] to convert the word sequence into the module program. However, there is no layout annotations in most real-world datasets (*e.g.*, VQAv2.0). In this case, following StackNMN [9], we make soft module selection with a differentiable stack structure. Please refer to their papers for more details.

We feed our output features from modules (cf. Table 1) into a softmax layer for the answer prediction. We use the cross entropy loss between our predicted answers and ground-truth answers to train our XNMs.

## 4. Experiments

### 4.1. CLEVR

**Settings.** The CLEVR dataset [12] is a synthetic diagnostic dataset that tests a range of visual reasoning abilities. In CLEVR, images are annotated with ground-truth object positions and labels, and questions are represented as functional programs that consists of 13 kinds of modules. Except the "Unique" module, which does not have actual operation, all the remaining 12 modules can correspond to our

Table 2: Comparisons between neural module networks on the CLEVR dataset. Top section: results of the official test set; Bottom section: results of the validation set (we can only evaluate our GT setting on the validation set since the annotations of the test set are not public [12]). The program option "scratch" means totally without program annotations, "supervised" means using trained end-to-end parser, and "GT" means using ground-truth programs. Our reasoning modules are composed with highly-reusable X modules, leading to a very small number of parameters. Using the ground-truth scene graphs and programs, we can achieve a perfect reasoning on all kinds of questions.

| Method | Program | #Modules | #Param. | Overall | Count | Compare Numbers | Exist | Query Attribute | Compare Attribute |
|---|---|---|---|---|---|---|---|---|---|
| Human [12] | - | - | - | 92.6 | 86.7 | 86.4 | 96.6 | 95.0 | 96.0 |
| N2NMN [10] | scratch | 12 | - | 69.0 | - | - | - | - | - |
| N2NMN [10] | supervised | 12 | - | 83.7 | - | - | - | - | - |
| PG+EE [13] | supervised | 39 | 40.4M | 96.9 | 92.7 | 98.7 | 97.1 | 98.1 | 98.9 |
| TbD-net [18] | supervised | 39 | 115M | **99.1** | **97.6** | **99.4** | **99.2** | **99.5** | **99.6** |
| StackNMN [9] | scratch | 9 | 7.32M | 93.0 | - | - | - | - | - |
| StackNMN [9] | supervised | 9 | 7.32M | 96.5 | - | - | - | - | - |
| XNM-Det | supervised | 12 | **0.55M** | 97.7 | 96.0 | 98.0 | 98.7 | 98.4 | 97.6 |
| NS-VQA [27] | supervised | 12 | - | 99.8 | 99.7 | 99.9 | 99.9 | 99.8 | 99.8 |
| XNM-Det | supervised | 12 | 0.55M | 97.8 | 96.0 | 98.1 | 98.6 | 98.7 | 97.8 |
| XNM-Det | GT | 12 | 0.55M | 97.9 | 96.2 | 98.1 | 98.8 | 98.7 | 97.8 |
| XNM-GT | supervised | 12 | 0.22M | 99.9 | 99.9 | 99.9 | 99.8 | 99.8 | 99.9 |
| XNM-GT | GT | 12 | **0.22M** | **100** | **100** | **100** | **100** | **100** | **100** |

modules in Table 1. CLEVR modules "Equal_attribute", "Equal_integer", "Greater_than" and "Less_than" have the same implementation as our `Compare`, but with different parameters. There are 4 attribute categories in CLEVR, so we set the number of mapping matrixes $K = 4$.

We reused the trained sequence-to-sequence program generator of [13, 18], which uses prefix-order traversal to convert the program trees to sequences. Note that their modules are bundled with input, e.g., they regard Filter[red] and Filter[green] as two different modules. This will cause serious sparseness in the real-world case. We used their program generator, but unpack the module and the input (e.g., Filter[red] and Filter[green] are the same module with different input query).

In the GT setting, we performed reasoning over the ground-truth scene graphs. In the Det setting, we built the scene graphs by detecting objects and using RoI features as node embeddings and the differences between detected coordinates as edge embeddings. Since CLEVR does not provide the bounding box or segmentation annotations of objects, it is hard to directly train an object detector. NS-VQA [27] trained a Mask R-CNN [7] for object segmentation by "hacking" the rendering process [12], which could perform very well due to the simplicity of visual scenes of CLEVR. However, as we expected to explore X modules in a noisier case, we chose the trained attention modules of TbD-net [18] as our object detector. Specifically, we enu-

merated all possible combinations of object attributes (e.g., red, cube, metal, large), and tried to find corresponding objects using their attention modules (e.g., intersection of the output mask of Attend[red], Attend[cube], Attend[metal] and Attend[large], and then regarded each clique as a single object). The detected results have some frequent mistakes, such as inaccurate position, wrongly merged nodes (two adjacent objects with the same attribute values are recognized as one). These detection noises allow us to test whether our XNMs are robust enough.

**Goals.** We expect to answer the following questions according to the CLEVR experiments: **Q1:** What is the upper bound of our X reasoning when both the vision and language perceptions are perfect? **Q2:** Are our XNMs robust for noisy detected scene graphs and parsed programs? **Q3:** What are the parameter and data efficiency, and the convergence speed of XNMs? **Q4:** How is the explainability of XNMs?

**Results.** Experimental results are listed in Table 2. **A1:** When using the ground-truth scene graphs and programs, we can achieve 100% accuracy, indicating an inspiring upper-bound of visual reasoning. By disentangling "high-level" reasoning from "low-level" perception and using XNMs, we may eventually conquer the visual reasoning challenge with the rapid development of visual recognition.

**A2:** With noisy detected scene graphs, we can still achieve a competitive 97.9% accuracy using the ground-

truth programs, indicating that our X reasoning are robust to different quality levels of scene graphs. When replacing the ground-truth programs with parsed programs, the accuracy drops by 0.1% in both GT and Det settings, which is caused by minor errors of the program parser.

**A3:** Due to the conciseness and high-reusability of X modules, our model requires significantly less parameters than existing models. Our GT setting only needs about 0.22M parameters, taking about 500MB memory with batch size of 128, while PG+EE [13] and TbD-net [18] bundle modules and inputs together, leading to a large number of modules and parameters.



Figure 3: Comparison of data efficiency and convergence speed.

To explore the data efficiency, we trained our model with a partial training set and evaluated on the complete validation set. Results are displayed in the left part of Figure 3. We can see that our model performs much better than other baselines when the training set is small. Especially, our GT setting can still achieve a 100% accuracy even with only 10% training data. The right part shows the accuracy at each training epoch. We can see that our X reasoning converges very fast.

**A4:** As our XNMs are attention-based, the reasoning process is totally transparent and we can easily show intermediate results. Figure 4 displays two examples of CLEVR. We can see all reasoning steps are clear and intuitive.

### 4.2. CLEVR-CoGenT

**Settings.** The CLEVR-CoGenT dataset is a benchmark to study the ability of models to recognize novel combinations of attributes at test-time, which is derived from CLEVR but has two different conditions: in Condition A all cubes are colored one of gray, blue, brown, or yellow, and all cylinders are one of red, green, purple, or cyan; in Condition B the color palettes are swapped. The model is trained using the training set of Condition A, and then is tested using Condition B to check whether it can generalize well to the novel attribute combinations. We train our model on the training set of Condition A, and report the accuracy

Table 3: Comparisons between NMNs on CLEVR-CoGenT. Top section: results of the test set; Bottom section: results of the validation set. Using the ground-truth scene graphs, our XNMs generalize very well and do not suffer from shortcuts at all.

| Method | Program | Condition A | Condition B |
|---|---|---|---|
| PG+EE [13] | supervised | 96.6 | 73.7 |
| TbD-net [18] | supervised | **98.8** | **75.4** |
| XNM-Det | supervised | 98.1 | 72.6 |
| NS-VQA [27] | supervised | 99.8 | 63.9 |
| XNM-Det | supervised | 98.2 | 72.1 |
| XNM-Det | GT | 98.3 | 72.2 |
| XNM-GT | supervised | 99.9 | 99.9 |
| XNM-GT | GT | **100** | **100** |

of both conditions.

**Goals. Q1:** Can our model perform well when meeting the novel attribute combinations? **Q2:** If not, what actually causes the reasoning shortcut?

**Results.** Results of CLEVR-CoGenT are displayed in Table 3. **A1:** When using the ground-truth scene graphs, our XNMs perform perfectly on both Condition A and Condition B. Novel combinations of attributes in Condition B do not cause the performance drop at all. However, when using the detected scene graphs, where node embeddings are RoI features that fuse all attribute values, our generalization results on Condition B drops to 72.1%, suffering from the dataset shortcut just like other existing models [13, 18].
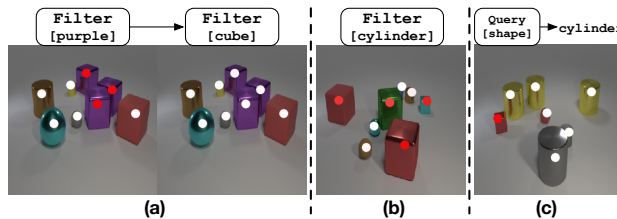


Figure 5: Failure cases of our Det setting on Condition B of CoGenT.

**A2:** Figure 5 shows some typical failure cases of our Det setting on Condition B. In case (a), our model cannot recognize purple cubes as "cube" because all cubes are colored one of gray, blue, brown, or yellow in the training data. Similarly, in case (b) and (c), whether an object is recognized as a "cube" or a "cylinder" by our model is actually determined by its color. However, in our GT setting, which is given the ground-truth visual labels, we can achieve a perfect performance. This gap reveals that the challenge of CLEVR-CoGenT mostly comes from the vision bias, rather than the reasoning shortcut.
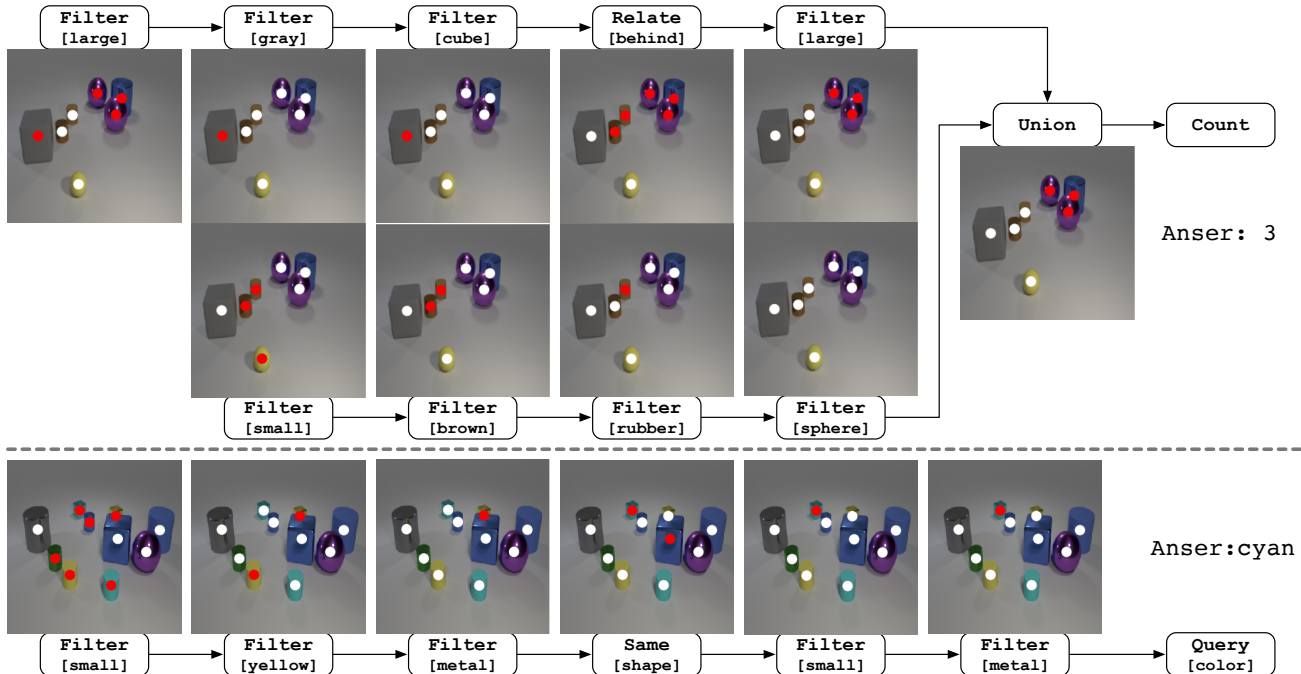
Figure 4: Reasoning visualizations of two CLEVR samples. Question 1: *What number of objects are either big objects that are behind the big gray block or tiny brown rubber balls?* Question 2: *The other small shiny thing that is the same shape as the tiny yellow shiny object is what color?* We plot a dot for each object and darker (red) dots indicate higher attention weights.

### 4.3. VQAv2.0

**Settings.** VQAv2.0 [6] is a real-world visual Q&A dataset which does not have annotations about scene graphs and module programs. We used the grounded visual features of [2] as node features, and concatenated node embeddings as edge features. We set $K = 1$ and fused the question embedding with our output feature for answer prediction. Following [2], we used softmax over objects for node attention computation.

**Goals.** We used VQAv2.0 to demonstrate the generality and robustness of our model in the practical case.

**Results.** We list the results in Table 4. We follow Stack-NMN [9] to build the module program in a stacked soft manner, but our model can achieve better performance as our reasoning over scene graphs is more powerful than their pixel-level operations.

Recall that [13, 18] are not applicable in open-vocabulary input, and [27] relies on the fixed label representation, so it is hard to apply them on practical datasets. In contrast, our XNMs are flexible enough for different cases.

## 5. Conclusions

In this paper, we proposed X neural modules (XNMs) that allows visual reasoning over scene graphs, represented by different detection qualities. Using the ground-truth

Table 4: Single-model results on VQAv2.0 validation set and test set. †: values reported in the original papers.

| Method | expert layout | validation(%) | test(%) |
|---|---|---|---|
| Up-Down [2] | no | 63.2† | 66.3 |
| N2NMN [10] | yes | - | 63.3† |
| StackNMN [9] | no | - | 64.1† |
| XNMs | no | **64.7** | **67.5** |

scene graphs and programs on CLEVR, we can achieve 100% accuracy with only 0.22M parameters. Compared to existing neural module networks, XNMs disentangle the "high-level" reasoning from the "low-level" visual perception, and allow us to pay more attention to teaching A.I. how to "think", regardless of what they "look". We believe that this is an inspiring direction towards explainable machine reasoning. Besides, our experimental results suggest that visual reasoning benefits a lot from high-quality scene graphs, revealing the practical significance of the scene graph research.

# References

[1] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *International Conference on Machine Learning*, 2016. 1

[2] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, 2018. 8

[3] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *CVPR*, 2016. 2

[4] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh. Vqa: Visual question answering. In *ICCV*, 2015. 1, 2

[5] L. Chen, H. Zhang, J. Xiao, X. He, S. Pu, and S.-F. Chang. Scene dynamics: Counterfactual critic multi-agent training for scene graph generation. *arXiv preprint arXiv:1812.02347*, 2018. 3

[6] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 1, 2, 3, 8

[7] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017. 6

[8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1

[9] R. Hu, J. Andreas, T. Darrell, and K. Saenko. Explainable neural computation via stack neural module networks. In *ECCV*, 2018. 2, 3, 5, 6, 8

[10] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. In *ICCV*, 2017. 2, 3, 6, 8

[11] D. A. Hudson and C. D. Manning. Compositional attention networks for machine reasoning. *ICLR*, 2018. 2

[12] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 1, 2, 3, 5, 6

[13] J. Johnson, B. Hariharan, L. van der Maaten, J. Hoffman, L. Fei-Fei, C. L. Zitnick, and R. B. Girshick. Inferring and executing programs for visual reasoning. In *ICCV*, 2017. 1, 2, 3, 6, 7, 8

[14] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei. Image retrieval using scene graphs. In *CVPR*, 2015. 3

[15] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 10

[16] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436, 2015. 1

[17] Y. Li, W. Ouyang, B. Zhou, J. Shi, C. Zhang, and X. Wang. Factorizable net: an efficient subgraph-based framework for scene graph generation. In *ECCV*, 2018. 3

[18] D. Mascharka, P. Tran, R. Soklaski, and A. Majumdar. Transparency by design: Closing the gap between performance and interpretability in visual reasoning. In *CVPR*, 2018. 2, 3, 6, 7, 8, 10

[19] M. L. Minsky. Logical versus analogical or symbolic versus connectionist or neat versus scruffy. *AI magazine*, 12(2):34, 1991. 1

[20] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a general conditioning layer. *AAAI*, 2018. 2

[21] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 1

[22] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. A simple neural network module for relational reasoning. In *NIPS*, 2017. 2

[23] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016. 1

[24] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014. 5

[25] D. Teney, L. Liu, and A. van den Hengel. Graph-structured representations for visual question answering. *arXiv preprint*, 2017. 3

[26] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei. Scene graph generation by iterative message passing. In *CVPR*, 2017. 3

[27] K. Yi, J. Wu, C. Gan, A. Torralba, P. Kohli, and J. B. Tenenbaum. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. *NIPS*, 2018. 2, 3, 6, 7, 8

[28] X. Yin and V. Ordonez. Obj2text: Generating visually descriptive language from object layouts. In *EMNLP*, 2017. 3

[29] R. Zellers, M. Yatskar, S. Thomson, and Y. Choi. Neural motifs: Scene graph parsing with global context. In *CVPR*, 2018. 2, 3

[30] Y. Zhang, J. Hare, and A. Prügel-Bennett. Learning to count objects in natural images for visual question answering. *arXiv preprint arXiv:1802.05766*, 2018. 4

# Appendix

## A. Implementation Details

In both of the GT and Det experiments on the CLEVR dataset, we set the dimension of the label embedding (*i.e.*, $d$) and the dimension of $\mathbf{h}$ in all output modules to 128. The classifier consists of a simple multi-layer perceptron that maps the 128-dimentional features to the number of possible answers (*i.e.*, 28), and a softmax layer. We used Adam [15] optimizer with an initial learning rate 0.001 to train our module parameters. We trained for 5 epochs for the GT setting and 10 epochs for the Det setting, and we reduced the learning rate to 0.0001 after the first epoch. Each epoch takes about one hour with an Nvidia 1080Ti graphic card.

The mapping matrices of the `Describe` module are implemented differently between the GT and Det setting. In the GT setting, as each object has four attribute values corresponding to four attribute categories (*i.e.*, color, shape, size, material), and our node embedding is the concatenation of attribute label embeddings, the dimension of node embeddings is $4d$, where $d$ is the dimension of label embeddings. We fix the order of label embedding concatenation as [color, shape, size, material], so we can extract node $i$'s color feature by $[\mathbf{I}_d; \mathbf{0}_d; \mathbf{0}_d; \mathbf{0}_d] \cdot \mathbf{v}_i$, where $\mathbf{I}_d, \mathbf{0}_d$ are $d \times d$ identity matrix and zero matrix respectively. So in the GT setting, our four mapping matrixes are defines as:

$$
\begin{aligned}
\mathbf{M}_1 &= [\mathbf{I}_d; \mathbf{0}_d; \mathbf{0}_d; \mathbf{0}_d], \\
\mathbf{M}_2 &= [\mathbf{0}_d; \mathbf{I}_d; \mathbf{0}_d; \mathbf{0}_d], \\
\mathbf{M}_3 &= [\mathbf{0}_d; \mathbf{0}_d; \mathbf{I}_d; \mathbf{0}_d], \\
\mathbf{M}_4 &= [\mathbf{0}_d; \mathbf{0}_d; \mathbf{0}_d; \mathbf{I}_d].
\end{aligned} \tag{9}
$$

In the Det setting, we regard $\mathbf{M}_k \in \mathbb{R}^{d \times d}$ as parameters and learn them automatically, which leads to an increase of the number of parameters (Det has 0.55M parameters while Gt only has 0.22M).

As for the attention functions in the CLEVR GT experiments, besides the label-space softmax which is mentioned in Eq. 5, we have also tried the sigmoid activation. Specifically, we fused multiple label vectors into one vector via a fully connected layer, and then applied the sigmoid function like Eq. 6 to separately compute attention weights of each node and edge. Using this attention strategy, we can still obtain 100% accuracy in the GT setting, demonstrating that our model is robust and flexible.

In the VQAv2.0 dataset, we selected the most frequent 3000 answers from the training set, and predicted the target answer from these candidates. We used an LSTM as the question encoder and fused the 1024-dimensional question embedding with the output feature from our module network for the answer classification. For the training, we used Adam optimizer with an initial learning rate 0.0008

and we set batch size as 256. The learning rate was decayed by half every 50000 training steps and the training lasted for 100 epochs.

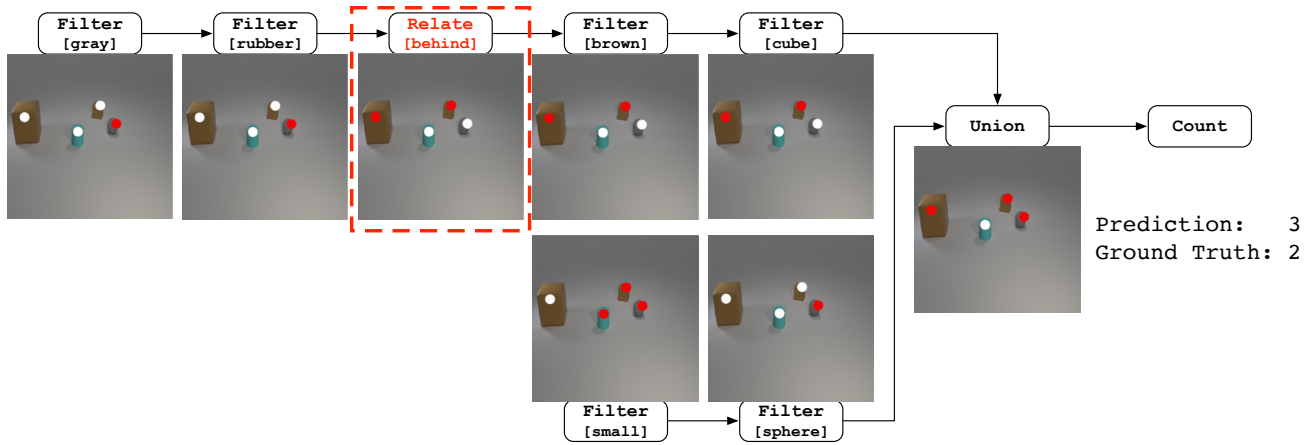## B. Failure Cases of the CLEVR Det Setting

We classify the failure cases in the CLEVR Det setting into three categories:

1. The coordinate detection is inaccurate (Figure 6).

2. Some objects are occluded (Figure 7).

3. The mask-based object division is inaccurate (Figure 8). Specifically, when we propose objects based on the generated mask from the "Attend" modules of [18], we may wrongly propose more or less objects due to the blurring boundaries.
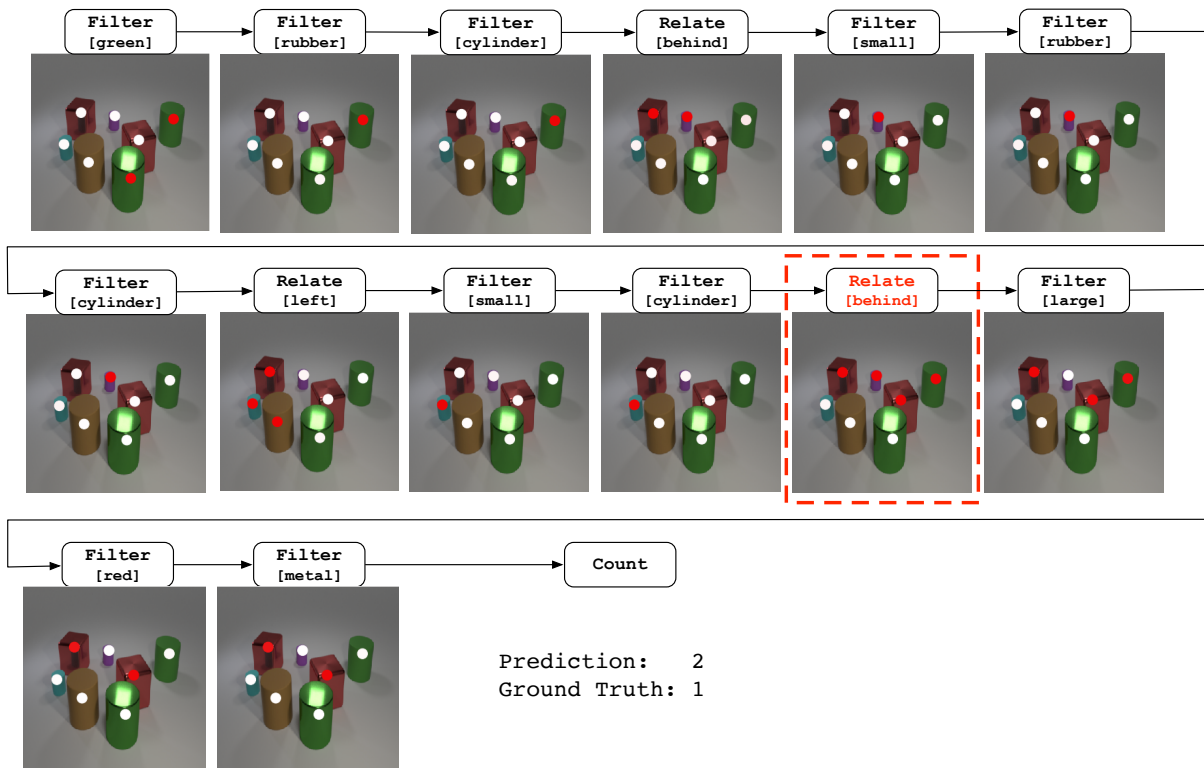
In all of Figure 6, 7, and 8, **we mark the reasoning steps that cause mistakes in red box**. We can see that using our X reasoning over scene graphs, we can easily track the reasoning process and diagnose where and why the model makes a mistake, which is an inspiring step towards the explainable AI.

## C. Case Study on the VQAv2.0

In the VQAv2.0 experiments, we predict a probability distribution over our modules at each step, and then feed the soft fusion of their outputs into next step. We set the reasoning length to 3 and force the last module to be `Describe`. We show a typical sample of the VQAv2.0 dataset in Figure 9. The top row is the modules with the most probability at each step, while the bottom row shows the results of `Relate` at Step 2. We can see that even though the question "what is the man wearing on his head" explicitly requires the relationship reasoning (*i.e.*, find the "man" first, and then move the focus to his head), our model scarcely relies on the results of the `Relate` module. Instead, it can directly focus on the "head" and give the correct prediction. We think this is due to the simplicity of questions, which is a shortcoming of the VQAv2.0 dataset.
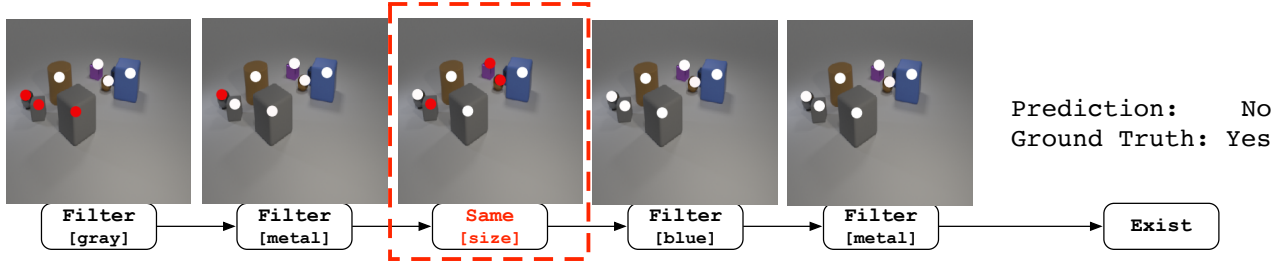
Question: What number of objects are tiny spheres or brown blocks behind the gray matte object ?
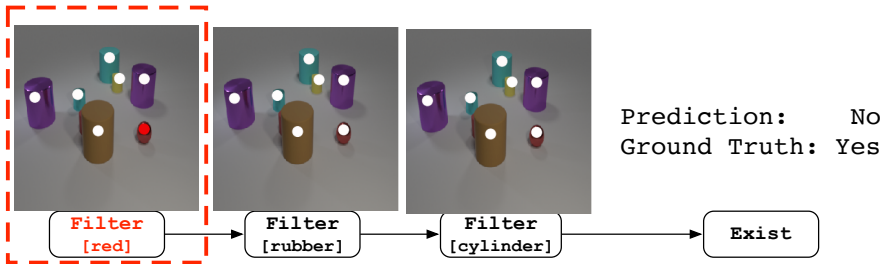


Question: There is a tiny cylinder left of the tiny matte cylinder that is behind the green rubber cylinder ; how many big red metal things are behind it ?

Figure 6: Failure cases caused by the inaccurate coordinate detection. Top case: the large brown cube is not behind the gray rubber object. Bottom case: a large red cube is wrongly recognized to be behind the tiny cylinder.

Prediction: No
Ground Truth: Yes

**Filter [gray]** → **Filter [metal]** → **Same [size]** → **Filter [blue]** → **Filter [metal]** → **Exist**

Question: Is there a blue metal object that has the same size as the gray metal object ?



Prediction: No
Ground Truth: Yes

**Filter [red]** → **Filter [rubber]** → **Filter [cylinder]** → **Exist**

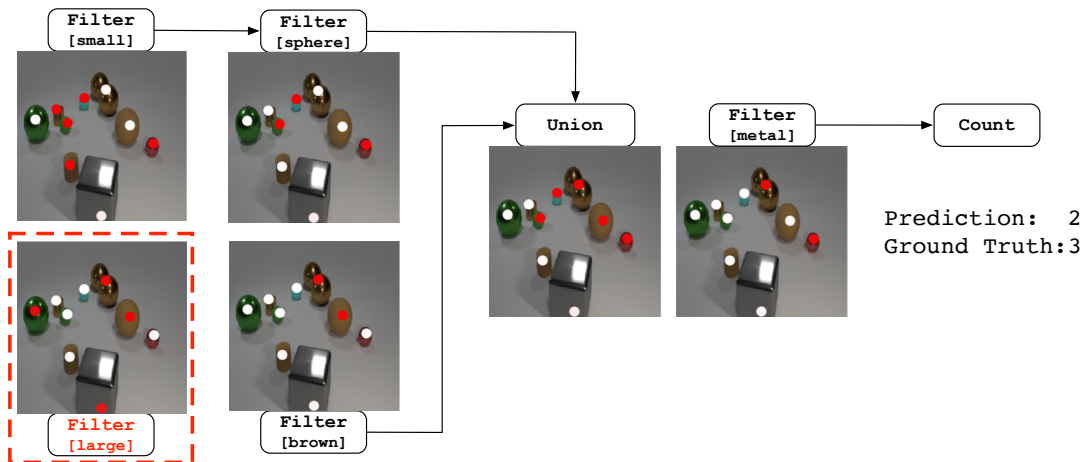Question: Are there any red rubber cylinders ?



Figure 7: Failure cases caused by occluded objects. We show the high-resolution images here, and we can see that 1) in the left image, there is a small blue cube behind the large blue cube occluded; 2) in the right image, there is a red cylinder behind the large brown cylinder occluded. These occluded objects do not have corresponding dots in the reasoning results, leading to a wrong prediction "No" while the actual answer is "Yes".

Question: There is a object that is both behind the cyan object and in front of the small metal cylinder ; what size is it ?



Question: What number of metal objects are large brown objects or tiny spheres ?

Figure 8: Failure cases caused by the inaccurate object division. Top case: two dots are assigned to the same object. Bottom case: two adjacent objects with the same attribute values (*i.e.*, large, brown, sphere, metal) are recognized as one object, which makes the predicted number less than the ground truth answer.

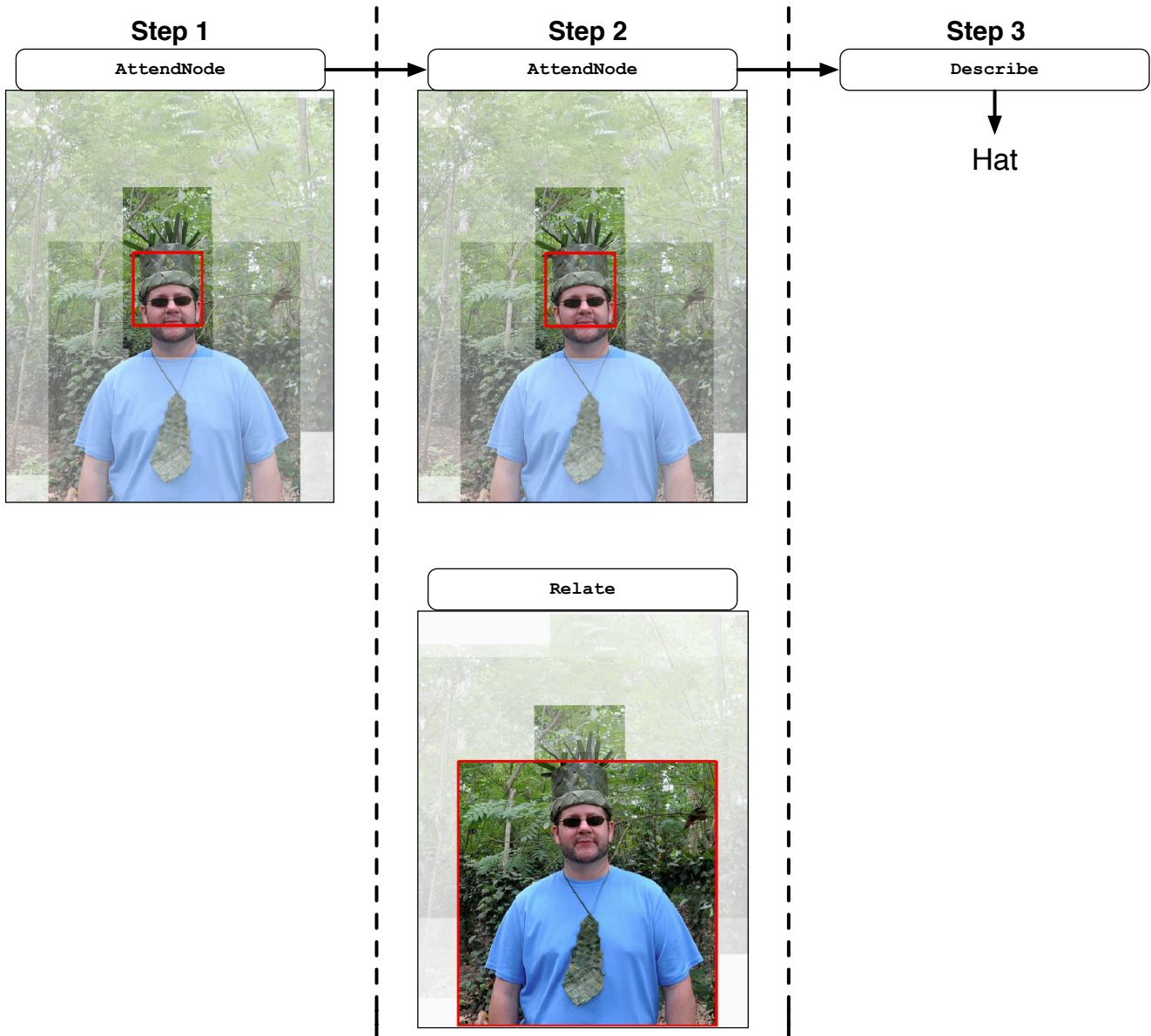Figure 9: A typical case of the VQAv2.0 dataset. The top row lists the modules with the maximum probability at each step. We can see that even the question "what is the man wearing on his head" explicitly requires the relationship understanding, the module Relate is still not necessary as the target region can be directly found. We argue the simplicity of question annotations is a major shortcoming of the VQAv2.0 dataset.