

CLEVR-Ref+: Diagnosing Visual Reasoning with Referring Expressions

Runtao Liu¹, Chenxi Liu^{2(✉)}, Yutong Bai³, Alan Yuille²

¹Peking University ²Johns Hopkins University ³Northwestern Polytechnical University

runtao219@gmail.com cxliu@jhu.edu ytongbai@gmail.com alan.1.yuille@gmail.com

Abstract

Referring object detection and referring image segmentation are important tasks that require joint understanding of visual information and natural language. Yet there has been evidence that current benchmark datasets suffer from bias, and current state-of-the-art models cannot be easily evaluated on their intermediate reasoning process. To address these issues and complement similar efforts in visual question answering, we build CLEVR-Ref+, a synthetic diagnostic dataset for referring expression comprehension. The precise locations and attributes of the objects are readily available, and the referring expressions are automatically associated with functional programs. The synthetic nature allows control over dataset bias (through sampling strategy), and the modular programs enable intermediate reasoning ground truth without human annotators.

In addition to evaluating several state-of-the-art models on CLEVR-Ref+, we also propose IEP-Ref, a module network approach that significantly outperforms other models on our dataset. In particular, we present two interesting and important findings using IEP-Ref: (1) the module trained to transform feature maps into segmentation masks can be attached to any intermediate module to reveal the entire reasoning process step-by-step; (2) even if all training data has at least one object referred, IEP-Ref can correctly predict no-foreground when presented with false-premise referring expressions. To the best of our knowledge, this is the first direct and quantitative proof that neural modules behave in the way they are intended.¹

1. Introduction

There has been significant research interest in the joint understanding of vision and natural language. While image captioning [17, 5, 25, 22] focuses on generating a sentence with image being the only input, visual question answering (VQA) [2, 6, 38] and referring expressions (REF) [24, 13] require comprehending both an image and a sentence, before generating an output. In this paper, we focus on refer-

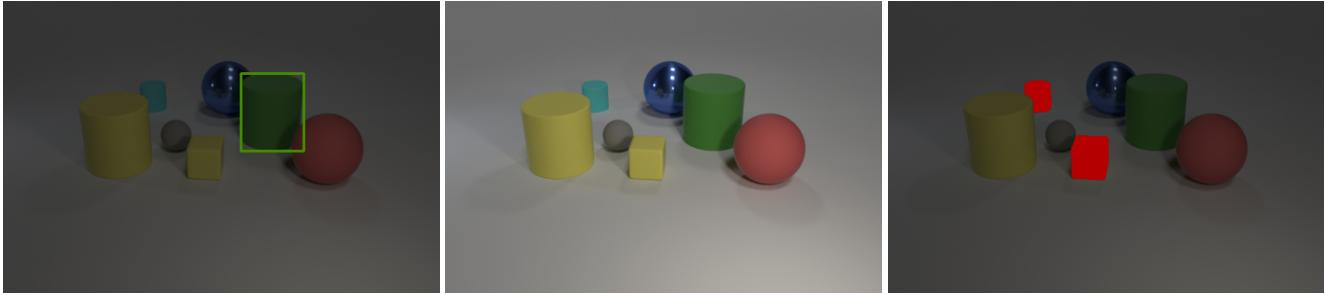
¹All data and code concerning CLEVR-Ref+ and IEP-Ref have been released at <https://cs.jhu.edu/~cxliu/2019/clevr-ref+>

ring expressions, which is to identify the particular objects (in the form of segmentation mask or bounding box) in a given scene from natural language.

In order to study referring expressions, various datasets have been proposed [24, 35, 18]. These are real-world images annotated by crowdsource workers. The advantage of these datasets is that they, to a certain extent, reflect the complexity and nuances of the real world. Yet inevitably, they also have limitations. First, they usually exhibit strong biases that may be exploited by the models [3]. Roughly speaking, this means simply selecting the salient foreground object (i.e., discarding the referring expression) will yield a much higher baseline than random. This casts doubts on the true level of understanding within current REF models. Second, evaluation can only be conducted on the final segmentation mask or bounding box, but not the intermediate step-by-step reasoning process. For example, for the referring expression “Woman to the left of the red suitcase”, a reasonable reasoning process should be first find all suitcases in the image, then identify the red one among them, finally segment the woman to its left. Clearly this requires significantly more high-quality annotations, which are currently unavailable and hard to collect.

To address these concerns and echo similar efforts in visual question answering (i.e., CLEVR [15]), we propose CLEVR-Ref+, a synthetic diagnostic dataset for referring expressions. The advantage of using a synthetic dataset is that we have full control over the scene, and dataset bias can be minimized by employing a uniform sampling strategy. Also, the referring expressions are now automatically annotated with the true underlying reasoning process, so a step-by-step analysis becomes much more plausible.

We make much effort in constructing CLEVR-Ref+ to make sure it is well adapted and applicable to the referring expression task. First, we turn the original questions in CLEVR into their corresponding referring expression format. Second, we change the output space from textual answers (in the form of a word) to referred objects (in the form of segmentation mask or bounding box). Third, we analyzed statistics from real-world REF datasets and found that there are some common types of referring expressions



The big thing(s) that are behind the second one of the big thing(s) from front and to the right of the first one of the large sphere(s) from left

Any other things that are the same size as the fifth one of the thing(s) from right

Figure 1: Examples from our CLEVR-Ref+ dataset. We use the same scenes as those provided in CLEVR [15]. Instead of asking questions about the scene, we ask the model to either return one bounding box (as illustrated on the left) or return a segmentation mask (could potentially be multiple objects; illustrated on the right) based on the given referring expression.

(e.g., “The second sphere from left”) that are not included in CLEVR templates. In our CLEVR-Ref+, we add support for these types of expressions to better match the variety of referring expressions used in real world.

We tested several state-of-the-art referring expression models on our CLEVR-Ref+ dataset. This includes both those designed for referring segmentation [21] and detection [36, 34]. In addition to evaluating the overall IoU and accuracy as previous datasets, we can now do a more detailed breakdown and analysis in terms of sub-categories. For example, we found that it is especially hard for the models to understand ordinality. This could point to important research directions in the future.

Besides diagnosing these existing models, we also propose IEP-Ref, a Neural Module Network [1] solution based on IEP [16]. Experiments show that the IEP-Ref model achieved excellent performance on CLEVR-Ref+ with its explicit, step-by-step functional program and module network execution engine, suggesting the importance of compositionality. Very interestingly, we found that the module trained on translating the last module output to segmentation mask is general, and can produce excellent human-interpretable segmentation masks when attached to intermediate module outputs, revealing the entire reasoning process. We believe ours is the first to show clean visualization of the visual reasoning process carried out by neural module networks, as opposed to gradient norms [16] or soft attention maps [27, 9].

In sum, our paper makes the following contributions:

- We construct CLEVR-Ref+, a synthetic diagnostic dataset for referring expression tasks that complements existing real-world datasets.
- We test and diagnose several state-of-the-art referring expression models on CLEVR-Ref+, including our proposed IEP-Ref that explicitly captures compositionality.

- The segmentation module trained in IEP-Ref can be trivially plugged in all intermediate steps in the module network to produce excellent segmentation masks that clearly reveal the network’s reasoning process.

2. Related Works

2.1. Referring Expressions

Referring expressions are sentences that refer to specific objects in an image. Understanding referring expressions has important applications in robotics and human-computer interaction. In recent years, many deep learning models have been developed for this task.

Several works focused on detection, i.e. returning one bounding box containing the referred object. [24, 13] adapted image captioning for this task by scoring each bounding box proposal with a generative captioning model. [32] learned the alignment between the description and image region by reconstructing the description using an attention mechanism. [35, 29] studied the importance of context for referring expressions. [23] used a discriminative comprehension model to improve referring expression generation. [36] showed additional gain by incorporating reinforcement learning. [11, 34] used learned parser and module networks to better match the structured semantics.

There are also works focusing on segmentation, i.e. returning the segmentation mask. [12] used FCN feature concatenated with LSTM feature to produce pixel-wise binary segmentation. [21] used a convolutional LSTM in addition to the language-only LSTM to facilitate propagation of intermediate segmentation beliefs. [20, 26] improved upon [21] by making more architectural improvements.

2.2. Dataset Bias and Diagnostic Datasets

In visual question answering, despite exciting models being proposed and accuracy on benchmark datasets being steadily improved, there has been serious concern over the

Table 1: Examples of converting questions to referring expressions.

Category	Question (CLEVR)	Referring Expression (CLEVR-Ref+)
Basic	How many cyan cubes are there?	The cyan cubes.
Spatial Relation	Are there any green cylinders to the left of the brown sphere?	The green cylinders to the left of the brown sphere.
AND Logic	How many green spheres are both in front of the red cylinder and left to the yellow cube?	The green spheres that are both in front of the red cylinder and left to the yellow cube.
OR Logic	Are there any cylinders that are either purple metal objects or small red matte things?	Cylinders that are either purple metal objects or small red matte things.
Same Relation	Are there any other things that have the same size as the red sphere?	The things/objects that have the same size as the red sphere.
Compare Integer	Are there more brown shiny objects behind the large rubber cylinder than gray blocks?	-
Comparison	Does the small ball have the same color as the small cylinder in front of the big sphere?	-

dataset bias problem [37, 7], meaning that models may be heavily exploiting the imbalanced distribution in the training/testing data. More recently, [3] showed that dataset bias also exists in referring expression datasets [24, 18, 35]. For example, [3] reported that the performance when discarding the referring expression and basing solely on the image is significantly higher than random. Ideally the dataset should be unbiased so that the performance faithfully reflect the model’s true level of understanding. But this is very hard to control when working with real-world images and human-annotated referring expressions.

A possible solution is to use synthetic datasets. Indeed this is the path taken by CLEVR [15], a diagnostic dataset for VQA. There, objects are placed on a 2D plane and only have a small number of choices in terms of shape, color, size, and material. The question-answer pairs are also synthesized using carefully designed templates. Together with a uniform sampling strategy, this design can mitigate dataset bias and reveal the model’s ability to understand compositionality. We construct our CLEVR-Ref+ dataset by repurposing CLEVR towards the referring expression task.

Several approaches now achieve near-perfect accuracy on CLEVR [16, 10, 30, 33, 27, 14, 9]. In addition to reporting the VQA accuracy, they typically try to interpret the visual reasoning process through visualization. However, the quality of these visualizations does not match the high VQA accuracy. We suspect the primary reason is that the domain these models are trained for (i.e. a textual answer) is different from the domain these models are diagnosed on (i.e. attention over the image). Fortunately, in referring expressions these two domains are very much interchangeable.

Note that CLEVR was also adapted towards referring expression in [9], but they focused on facilitating VQA, instead of introducing extensions (Section 3.3), evaluating state-of-the-art models (Section 4.1), and directly facilitating the diagnosis of visual reasoning (Section 4.3).

3. The CLEVR-Ref+ Dataset

CLEVR-Ref+ uses the exact same scenes as CLEVR (70K images in train set, 15K images in validation and test set), and every image is associated with 10 referring expressions. Since CLEVR is a VQA dataset, we began by changing the questions to referring expressions (Section 3.1), and the answers to referred objects (Section 3.2). We then made important additions to the set of modules (Section 3.3) as well as necessary changes to the sampling procedure (Section 3.4). Finally, we made the distinction whether more than one object is being referred (Section 3.5).

3.1. From Question to Referring Expression

Templates are provided in CLEVR so that questions and the functional programs associated with them can be generated at the same time. We notice that in many cases, part of the question is indeed a referring expression, as we need to first identify objects of interest before asking about their property (e.g. color or number). In Table 1 we provide examples of how we change question templates into their corresponding referring expression templates, usually by selecting a subset. The associated functional programs are also adjusted accordingly. For example, for “How many” questions, we simply remove the Count module at the end.

The original categories “Compare Integer” and “Comparison” were about comparing properties of two groups of referred objects, so they do not contribute additional referring expression patterns. Therefore they are not included in the templates for CLEVR-Ref+.

3.2. From Answer to Referred Objects

In referring expressions, the output is no longer a textual answer, but a bounding box or segmentation mask.

Since we know the exact 3D locations and properties of objects in the scene, we can follow the ground truth func-

Table 2: Frequent category and words in RefCOCO+ [35].

Category	Example words	Frequency
object	shirt,head,chair,hat,pizza	63.66%
human	man,woman,guy,girl,person	42.54%
color	white,black,blue,red,green	38.76%
spatial	back,next,behind,near,up	23.86%
animal	zebra,elephant,horse,bear	15.36%
attribute	big,striped,small,plaid,long	10.55%
action	standing,holding,looking	10.34%
ordinal	closest,furthest,first,third	5.797%
compare	smaller,tallest,shorter,older	5.247%
visible	fully visible,barely seen	4.639%

tional program associated with the referring expression to identify which objects are being referred. In fact we can do this not only at the end (also available in real-world datasets), but also at every intermediate step (not available in real-world datasets). This will become useful later when we do step-by-step inspection and evaluation of the visual reasoning process.

After finding the referred objects, we project them back to the image plane to get the ground truth bounding box and segmentation mask. This automatic annotation was done through rendering with the software Blender. For occluded objects, only the visible part is treated as ground truth.

3.3. Module Additions

We hope the referring expressions that we generate are representative of those used in the real world. However, since the task is no longer the same, we suspect that there may be some frequent referring patterns missing in the templates directly inherited from CLEVR. To this end, we analyzed statistics from a real-world referring expression dataset, RefCOCO+ [35], as shown in Table 2.

We began by sorting the words in these referring expressions by their frequency. Then, starting with the most frequent word, we empirically cluster these words into categories. Not surprisingly, nouns that represent object or human are the most common. However, going down the list, we found that the “ordinal” (e.g. “The second woman from left”) and “visible” (e.g. “The barely seen backpack”) categories recall more than 10% of all sentences, but are not included in the existing templates. Moreover, it is indeed possible to define them using a computer program, because there is no ambiguity in meaning. We add these two new modules into the CLEVR-Ref+ function catalog.

In a functional program, these two modules may be inserted whenever color, material, size, or shape is being described. As an example, “the red sphere” may be equivalently described as “the third sphere from left” or “the partially visible red object”. In our dataset, we define an object

to be *partially visible* if foreground objects’ mask occupies more than 20% of its bounding box area. For an object to be *fully visible*, this value must be exactly 0. We do not describe visibility when there is an ambiguous case (i.e. this value is between 0 and 0.2) in the scene.

3.4. Generation Procedure

Generating a referring expression for a scene is conceptually simple and intuitive. The process may be summarized as the following few steps:

1. Randomly choose a referring expression family².
2. Randomly choose a text template from this family.
3. Follow the functional program and select random values when encountering template parameters³.
4. Reject when certain criteria fail, that is, the sampled referring expression is inappropriate for the given scene; return when the entire functional program follows through.

We largely follow the generation procedure of CLEVR, with a few important changes:

- To balance the number of referring expressions across different categories (those listed in Table 1), we double the probability of being sampled in categories with a small number of referring expression families.
- When describing the attributes for a set of objects, we do not use `Ordinal` and `Visible` at the same time. This is because referring an object as “The second partially visible object from left” seems too peculiar and rare, and there usually exists more natural alternatives.
- Originally when describing the attributes for a set of objects, four fair coins were flipped to determine whether color, material, size, shape will be included. As a result, usually multiple attributes are selected, and a very small number of objects survive these filters. We empirically found that this makes it quite easy for the system to select the correct object simply from the attributes that directly describe the target object(s).

To remedy this, we first enumerate all possible combinations of these attributes, and calculate how many objects will survive for each possibility. We then uniformly sample from these possible number of survivors, before doing another uniform sampling to find the combination of attributes. This will ensure a larger variance in terms of number of objects after each set of filtering, and prevent near-degenerate solutions.

- At the end of the functional program, we verify if at least one object is being referred; reject otherwise.

²A referring expression family contains a template for constructing functional programs and several text templates that provide multiple ways of expressing these programs in natural language.

³For instance, left/right/front/behind; big/small; metal/rubber.

Table 3: Referring object detection and referring image segmentation results on CLEVR-Ref+. We evaluated three existing models, as well as IEP-Ref which we adapted from its VQA counterpart.

	Basic 0-Relate	Spatial Relation			Logic		Same	Accuracy	IoU
		1-Relate	2-Relate	3-Relate	AND	OR			
SLR [36]	0.627	0.569	0.570	0.584	0.594	0.701	0.444	0.577	-
MAttNet [34]	0.566	0.623	0.634	0.624	0.723	0.737	0.454	0.609	-
RMI [21]	0.822	0.713	0.736	0.715	0.585	0.679	0.251	-	0.561
IEP-Ref (GT)	0.928	0.895	0.908	0.908	0.879	0.881	0.647	-	0.816
IEP-Ref (700K prog.)	0.920	0.884	0.902	0.898	0.860	0.869	0.636	-	0.806
IEP-Ref (18K prog.)	0.907	0.858	0.874	0.862	0.829	0.847	0.605	-	0.782
IEP-Ref (9K prog.)	0.910	0.858	0.847	0.811	0.778	0.791	0.626	-	0.760

3.5. Multi-Object and Single-Object Referring

As explained in Section 3.4, each referring expression in CLEVR-Ref+ may refer to one or more objects in the scene. We believe this is the more general setting, and models should have the flexibility to handle various number of objects being referred. This is already handled and supported by referring image segmentation systems. However, we notice that detection based systems are usually designed to return a single object instead of multiple objects, presumably because this was how the detection datasets [24, 35] were created. As a result, for detection based methods, we evaluate on the subset of CLEVR-Ref+ where only a single object is referred. This subset contains a total of 222,569 referring expressions (32% of the entire dataset).

4. Experiments

4.1. Models and Implementation Details

In all models we resize the input image to 320×320 to set up a fair comparison. Publicly available code for these models are used with minimum change to adapt to our CLEVR-Ref+ dataset. The following referring expression models are studied and tested:

Speaker-Listener-Reinforcer (SLR) [36] This is a **detection** model that includes a generative model (speaker), a discriminative model (listener), as well as a reinforcement learning component that makes further improvement. Before training the main model, the visual-language similarity model needs to be trained first. We use Adam optimizer [19], learning rate 4e-4, batch size 32 for both the visual-language similarity model and the main model.

MAttNet [34] This is also a **detection** model, that uses three modular networks to capture the subject, location, and relationship features respectively. A soft attention mechanism is used to return the overall score of a candidate region. We use learning rate 4e-4 and batch size 15.

Recurrent Multimodal Interaction (RMI) [21] This is a **segmentation** model. In addition to concatenating the refer-

ring expression LSTM embedding with the image features, RMI also used a convolutional LSTM to facilitate propagation of segmentation beliefs when reading in the referring expression word-by-word. We use Adam optimizer, learning rate 2.5e-4, batch size 3, and weight decay 5e-4.

IEP-Ref This is a **segmentation** model that we adapt from IEP [16], which was designed for VQA. The idea is to use a LSTM program generator to translate the referring expression into a structured series of modules, each of which is parameterized by a small CNN. By executing this dynamically constructed neural network (with a special Segment module at the end; see supplementary material for its architecture), IEP-Ref imitates the underlying visual reasoning process. For input visual features, we use the last layer of the conv4 stage of ResNet101 [8] pre-trained on ImageNet [4], which is of size $1024 \times 20 \times 20$. Following [16], this part is not finetuned. We tried three settings that use 9K/18K/700K ground truth programs to train the LSTM program generator (Adam optimizer, learning rate 5e-4, batch size 64; 20,000 iterations for the 9K setting, 32,000 iterations for the 18K and 700K setting). The accuracies of the predicted programs are 0.873, 0.971, 0.993 respectively. For the fourth setting, we simply use the ground truth program⁴. The execution engine is trained for 30 epochs using learning rate 1e-4 and Adam optimizer.

4.2. Results and Analysis

4.2.1 Overall Evaluation

The experimental results are summarized in Table 3. Detection models are evaluated by accuracy (i.e. whether the prediction selects the correct bounding box among given candidates), where MAttNet performs favorably against SLR. Segmentation models are evaluated by Intersection over Union (IoU), where IEP-Ref performs significantly better than RMI. This suggests the importance to model compositionality within the referring expression. We now present a more detailed analysis of various aspects.

⁴This is our default IEP-Ref setting unless otherwise specified.

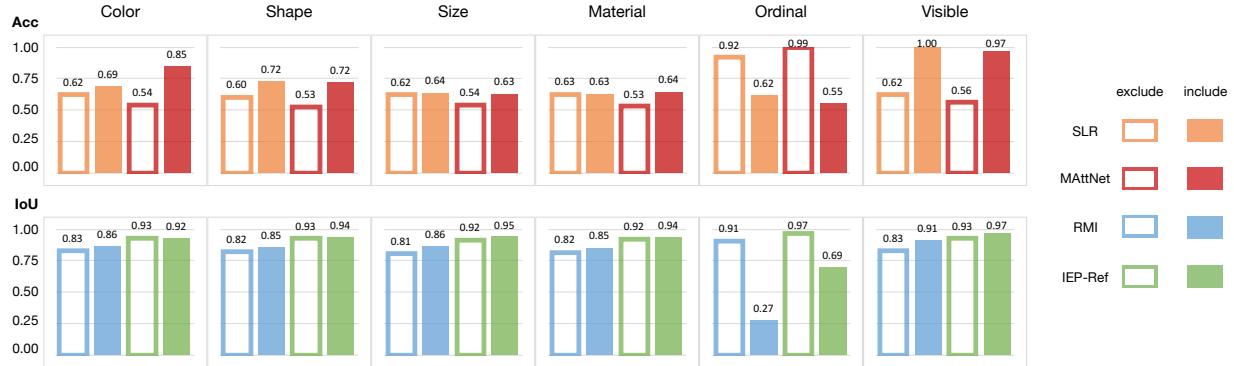


Figure 2: Analyzing the basic referring ability of different models. “Include” means the average performance if a module is involved in the referring process. “Exclude” means otherwise. As a result, high “exclude” and low “include” performance suggests that this module is more challenging to learn, and vice versa.

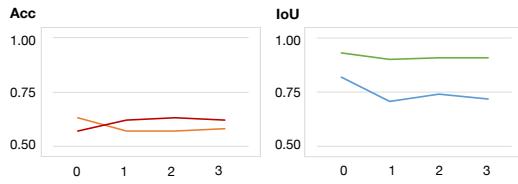


Figure 3: Analyzing the spatial reasoning ability of different models. Horizontal axis is the number of spatial relations.

4.2.2 Basic Referring Ability

Here we start with the easiest form: referring by direct description of object attributes (e.g., “The big blue sphere”). Concretely, this corresponds to the “0-Relate” subset.

In CLEVR-Ref+, there are totally 6 types of attributes that may help us locate specific objects: color, size, shape, material, ordinality, and visibility. In Figure 2 we show the average detection accuracy/segmentation IoU of various methods on “0-Relate” referring expressions that either contain or not contain a specific type of module.

Among detection models, we found that accuracy is higher when the referring expression contains descriptions of color, shape, and visibility. A reasonable conjecture is that these concepts are easier to learn compared with the others. However, for segmentation, the performance gaps between “exclude” and “include” are not as significant.

Though it is unclear which concept is the easiest to learn, there seems little dispute that ordinality is the hardest. In particular, for RMI, IoU is 0.91 if the expression does not require ordinality and 0.27 when it does. Other models do not suffer as much, but also experience significant drops. We suspect this is because ordinality requires the global context, whereas the others are local properties.

4.2.3 Spatial Reasoning Ability

Other than directly describing the attributes, it is also common to refer to an object by its spatial location. Here we diagnose whether referring expression models can understand

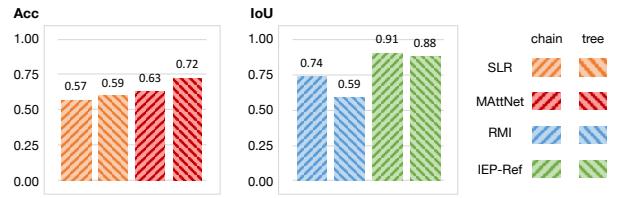


Figure 4: Effect of reasoning topology (Chain vs. Tree) on referring detection or segmentation performance.

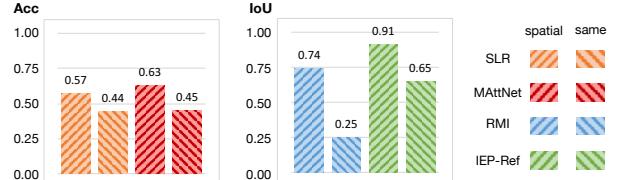


Figure 5: Effect of relation type (Spatial vs. Same) on referring detection or segmentation performance.

(potentially multiple steps of) relative spatial relationship, for example “The object that is left to the red cube”. In Table 3, this corresponds to the “{0, 1, 2, 3}-Relate” columns. Results are shown in Figure 3.

In general, we observe a small drop when referring expressions start to include spatial reasoning. However, there does not seem to be significant difference among referring expressions that require 1, 2, 3 steps of spatial reasoning. This seems to suggest that once the model has grasped spatial reasoning, there is little trouble in successfully applying it multiple times.

4.2.4 Different Reasoning Topologies

There are two referring expression topologies in CLEVR-Ref+: chain-structured and tree-structured. Intuitively, a chain structure has a single reasoning path to follow, whereas a tree structure requires following two such paths before merging. In Figure 4 we compare performance on referring expressions with two sequential spatial relation-

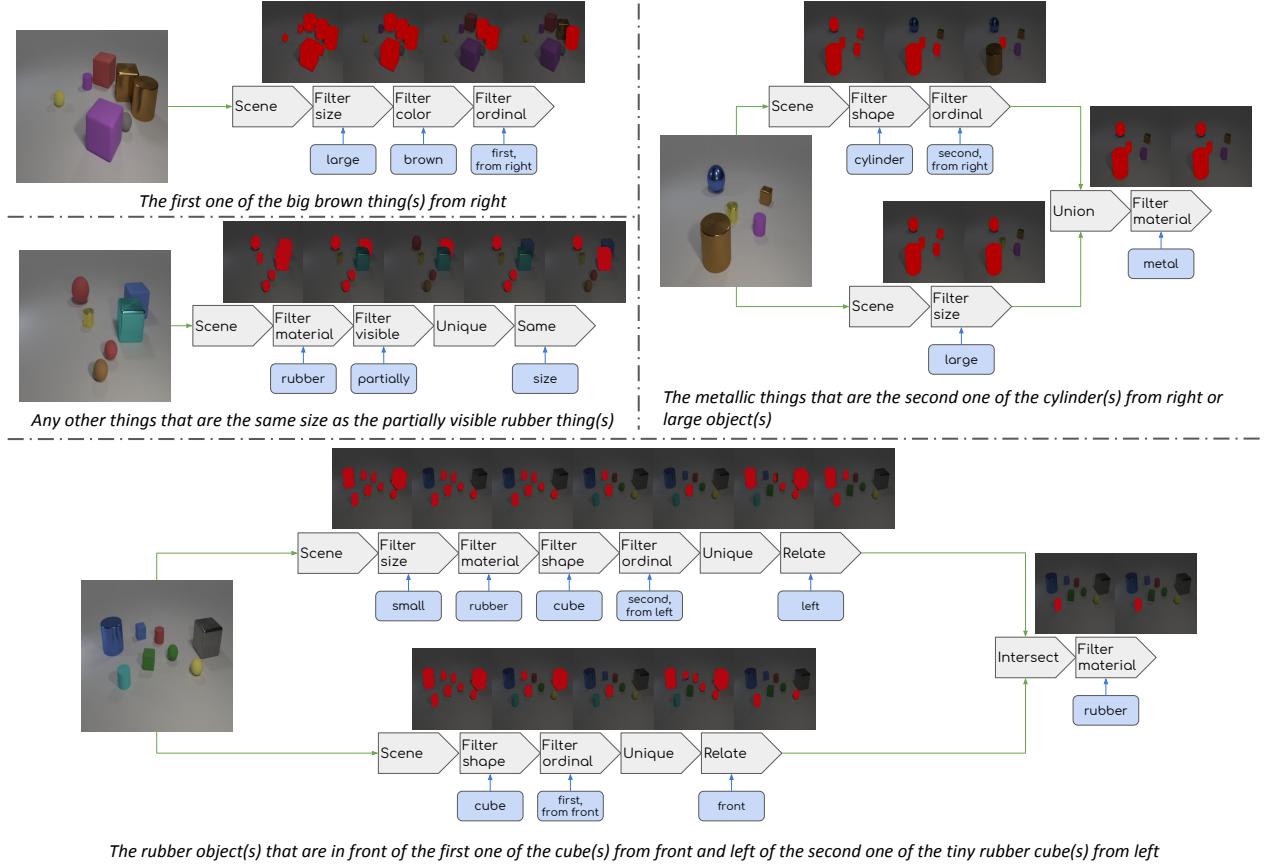


Figure 6: Four examples (two chain structures, two tree structures) of step-by-step inspection of IEP-Ref visual reasoning.

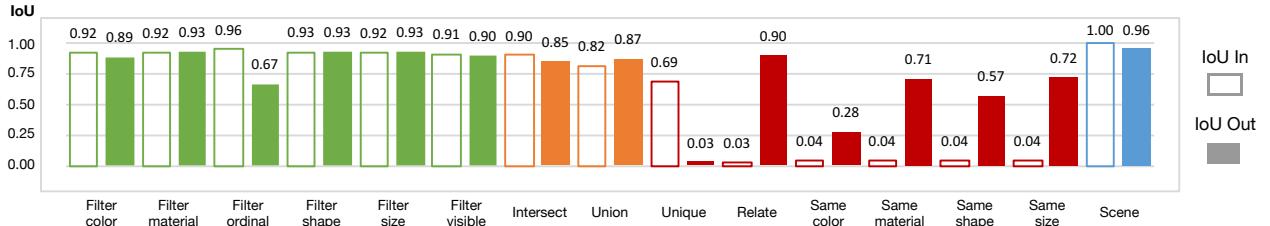


Figure 7: Average IoU going into/out of each IEP-Ref module on CLEVR-Ref+ validation set. Note that here IoU is not only computed at the end, but also all intermediate steps. This shows that IoU remains high throughout visual reasoning. The large differences in modules marked in dark red are discussed in text.

ships vs. one on each branch joined with AND. These two templates have roughly the same length and complexity, so the comparison focuses on topology.

Though not consistent among the four models, tree-structured referring expressions are generally harder than chain-structured ones. This agrees with the findings in [15].

4.2.5 Different Relation Types

There are two kinds of relationships in CLEVR-Ref+. One is spatial relationship that includes phrases like “left of”, “right of”, “in front of”, “behind” (discussed in Sec-

tion 4.2.3). The other is same-attribute relationship that requires recognizing and memorizing particular attributes of another object, e.g. “The large block(s) that have the same color as the metal sphere”.

In Figure 5 we study whether the relation type will make a difference in performance. We compare the “2-Relate” column with the “Same” column in Table 3, again because they have roughly the same length and complexity. All models perform much worse on the same-attribute relationship type, suggesting that this is a hard concept to grasp. Similar to ordinality, same-attribute requires global context.

4.3. Step-By-Step Inspection of Visual Reasoning

All the results discussed in Section 4.2 are about the endpoint of the visual reasoning process. We argue that in order to trust the predictions made by the referring expression system, it is also important to make sure that the intermediate reasoning steps make sense. CLEVR-Ref+ is suitable because: (1) the semantics of the referring expressions is modularized, and (2) the referring ground truth at all intermediate steps can be obtained automatically (i.e. no human annotators needed).

In training our IEP-Ref model, there is always a Segment module at the end, transforming the 128-channel feature map into a 1-channel segmentation mask. When testing, we simply attach the trained Segment module to the output of all intermediate modules. This is possible because all modules have the same number of input channels and output channels (128). This technique would not help in the VQA setting, because there the ending modules (e.g. Count, Equal) discard the spatial dimensions needed for visualization.

We found that this technique works quite well. In Figure 6 we provide four qualitative examples with various topologies and modules. We notice that all modules are performing their intended functionality, except the Unique module⁵. Yet after one more module, the segmentation mask becomes normal again. The quantitative analysis in Figure 7 confirms this observation: on average, IoU drops by 0.66 after each Unique module; but IoU significantly increases after each Same or Relate module, and these are the only modules that may come after Unique according to the templates. We conjecture that the network has learned some mechanism to treat Unique as the “preprocessing” step of the Same and Relate functionalities.

4.4. False-Premise Referring Expressions

In reality, referring expression systems may face all kinds of textual input, and not all of them will make sense. When presented with a referring expression that makes false assumptions (e.g. “The red sphere” when there is no sphere in the scene), the system should follow through as much as it can, and be robust enough to return zero foreground at the end. We test IEP-Ref’s ability to deal with these false-premise referring expressions (c.f. [31]). Note that no such expressions appear during training.

We generate 10,000 referring expressions that refer to zero object at the end. Qualitatively (see Figure 8), it is reassuring to see that intermediate modules are correctly doing their jobs, and a no-foreground prediction is made at the final step. Quantitatively, IEP-Ref predicts 0 foreground

⁵It is supposed to simply carry over the previously referred object, yet from what we observe, its behavior is most similar to selecting the complement of the previously referred object, though this is far from consistent.

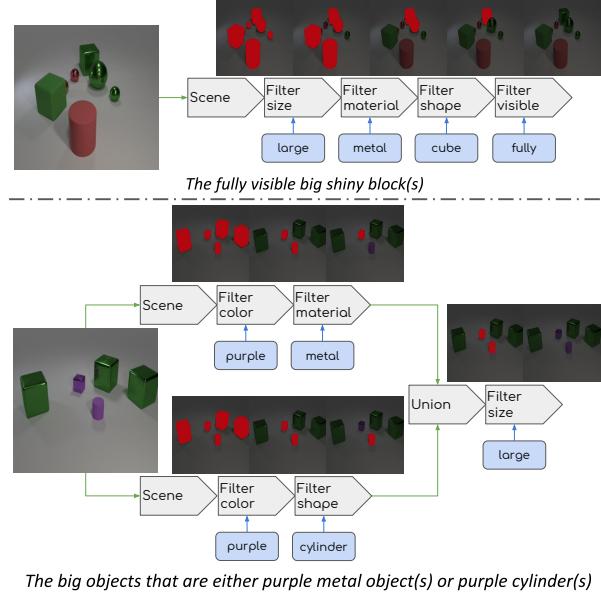


Figure 8: Our IEP-Ref model can correctly handle false-premise referring expressions even if they do not appear during training.

pixel more than 1/4 of the time, and ≤ 8 foreground pixels more than 1/3 of the time.

5. Conclusion

In this paper, we build the CLEVR-Ref+ dataset to complement existing ones for referring expressions. By choosing a synthetic setup, the advantage is that dataset bias can be minimized, and the ground truth visual reasoning process is readily available. We evaluated several state-of-the-art referring object detection and referring image segmentation models on CLEVR-Ref+. In addition, we propose the IEP-Ref model, which uses a module network approach and outperforms competing methods by a large margin. Detailed analysis are conducted to identify the strengths and weaknesses of these models. In particular, we found that ordinality and the same-attribute relationship seem to be the most difficult concepts to grasp.

Besides the correctness of the final segmentation mask, the correctness of the reasoning process is also important. We discovered that IEP-Ref provides an easy and natural way of revealing this process: simply attach the Segment module to each intermediate step. Our quantitative evaluation shows a high IoU at intermediate steps as well, proving that the neural modules have indeed learned the job they are supposed to do. Another evidence is that IEP-Ref can correctly handle false-premise referring expressions.

Going forward, we are interested to see whether these findings will transfer and inspire better models on real data.

Acknowledgments This research is support by NSF award CCF-1317376 and ONR N00014-12-1-0883.

References

- [1] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *CVPR*, pages 39–48. IEEE Computer Society, 2016. [2](#)
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: visual question answering. In *ICCV*, pages 2425–2433. IEEE Computer Society, 2015. [1](#)
- [3] V. Cirik, L. Morency, and T. Berg-Kirkpatrick. Visual referring expression recognition: What do systems actually learn? In *NAACL-HLT (2)*, pages 781–787. Association for Computational Linguistics, 2018. [1, 3](#)
- [4] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE Computer Society, 2009. [5](#)
- [5] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, T. Darrell, and K. Saenko. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, pages 2625–2634. IEEE Computer Society, 2015. [1](#)
- [6] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question. In *NIPS*, pages 2296–2304, 2015. [1](#)
- [7] 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*, pages 6325–6334. IEEE Computer Society, 2017. [3](#)
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778. IEEE Computer Society, 2016. [5, 11](#)
- [9] R. Hu, J. Andreas, T. Darrell, and K. Saenko. Explainable neural computation via stack neural module networks. In *ECCV (7)*, volume 11211 of *Lecture Notes in Computer Science*, pages 55–71. Springer, 2018. [2, 3](#)
- [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*, pages 804–813. IEEE Computer Society, 2017. [3](#)
- [11] R. Hu, M. Rohrbach, J. Andreas, T. Darrell, and K. Saenko. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, pages 4418–4427. IEEE Computer Society, 2017. [2](#)
- [12] R. Hu, M. Rohrbach, and T. Darrell. Segmentation from natural language expressions. In *ECCV (1)*, volume 9905 of *Lecture Notes in Computer Science*, pages 108–124. Springer, 2016. [2](#)
- [13] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell. Natural language object retrieval. In *CVPR*, pages 4555–4564. IEEE Computer Society, 2016. [1, 2](#)
- [14] D. A. Hudson and C. D. Manning. Compositional attention networks for machine reasoning. *CoRR*, abs/1803.03067, 2018. [3](#)
- [15] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. B. Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, pages 1988–1997. IEEE Computer Society, 2017. [1, 2, 3, 7, 12](#)
- [16] 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*, pages 3008–3017. IEEE Computer Society, 2017. [2, 3, 5, 11](#)
- [17] A. Karpathy and F. Li. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, pages 3128–3137. IEEE Computer Society, 2015. [1](#)
- [18] S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, pages 787–798. ACL, 2014. [1, 3](#)
- [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. [5](#)
- [20] R. Li, K. Li, Y.-C. Kuo, M. Shu, X. Qi, X. Shen, and J. Jia. Referring image segmentation via recurrent refinement networks. In *CVPR*, pages 5745–5753. IEEE Computer Society, 2018. [2](#)
- [21] C. Liu, Z. Lin, X. Shen, J. Yang, X. Lu, and A. L. Yuille. Recurrent multimodal interaction for referring image segmentation. In *ICCV*, pages 1280–1289. IEEE Computer Society, 2017. [2, 5](#)
- [22] C. Liu, J. Mao, F. Sha, and A. L. Yuille. Attention correctness in neural image captioning. In *AAAI*, pages 4176–4182. AAAI Press, 2017. [1](#)
- [23] R. Luo and G. Shakhnarovich. Comprehension-guided referring expressions. In *CVPR*, pages 3125–3134. IEEE Computer Society, 2017. [2](#)
- [24] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, pages 11–20. IEEE Computer Society, 2016. [1, 2, 3, 5](#)
- [25] J. Mao, W. Xu, Y. Yang, J. Wang, and A. L. Yuille. Deep captioning with multimodal recurrent neural networks (m-rnn). *CoRR*, abs/1412.6632, 2014. [1](#)
- [26] E. Margffoy-Tuay, J. C. Pérez, E. Botero, and P. Arbeláez. Dynamic multimodal instance segmentation guided by natural language queries. In *ECCV (11)*, volume 11215 of *Lecture Notes in Computer Science*, pages 656–672. Springer, 2018. [2](#)
- [27] D. Mascharka, P. Tran, R. Soklaski, and A. Majumdar. Transparency by design: Closing the gap between performance and interpretability in visual reasoning. *CoRR*, abs/1803.05268, 2018. [2, 3](#)
- [28] I. Misra, R. B. Girshick, R. Fergus, M. Hebert, A. Gupta, and L. van der Maaten. Learning by asking questions. In *CVPR*, pages 11–20. IEEE Computer Society, 2018. [12](#)
- [29] V. K. Nagaraja, V. I. Morariu, and L. S. Davis. Modeling context between objects for referring expression understanding. In *ECCV (4)*, volume 9908 of *Lecture Notes in Computer Science*, pages 792–807. Springer, 2016. [2](#)
- [30] E. Perez, F. Strub, H. de Vries, V. Dumoulin, and A. C. Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI*, pages 3942–3951. AAAI Press, 2018. [3](#)
- [31] A. Ray, G. Christie, M. Bansal, D. Batra, and D. Parikh. Question relevance in VQA: identifying non-visual and

- false-premise questions. In *EMNLP*, pages 919–924. The Association for Computational Linguistics, 2016. 8
- [32] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele. Grounding of textual phrases in images by reconstruction. In *ECCV(1)*, volume 9905 of *Lecture Notes in Computer Science*, pages 817–834. Springer, 2016. 2
- [33] A. Santoro, D. Raposo, D. G. T. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. A simple neural network module for relational reasoning. In *NIPS*, pages 4974–4983, 2017. 3
- [34] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg. Mattnet: Modular attention network for referring expression comprehension. In *CVPR*. IEEE Computer Society, 2018. 2, 5
- [35] L. Yu, P. Poirson, S. Yang, A. C. Berg, and T. L. Berg. Modeling context in referring expressions. In *ECCV (2)*, volume 9906 of *Lecture Notes in Computer Science*, pages 69–85. Springer, 2016. 1, 2, 3, 4, 5
- [36] L. Yu, H. Tan, M. Bansal, and T. L. Berg. A joint speaker-listener-reinforcer model for referring expressions. In *CVPR*, pages 3521–3529. IEEE Computer Society, 2017. 2, 5
- [37] P. Zhang, Y. Goyal, D. Summers-Stay, D. Batra, and D. Parikh. Yin and yang: Balancing and answering binary visual questions. In *CVPR*, pages 5014–5022. IEEE Computer Society, 2016. 3
- [38] Y. Zhu, O. Groth, M. S. Bernstein, and L. Fei-Fei. Visual7w: Grounded question answering in images. In *CVPR*, pages 4995–5004. IEEE Computer Society, 2016. 1

Supplementary Material

In this supplementary material, we begin by providing network architecture details of IEP-Ref to supplement Section 4.1 of the main paper. We then provide more analysis of the four models’ performance on CLEVR-Ref+, to supplement Section 4.2 of the main paper. Finally, we show more qualitative examples (referring expression and ground truth box/mask) from CLEVR-Ref+.

A. Network Architectures in IEP-Ref

In Figure 7 of the main paper, we listed all modules used in our IEP-Ref model (except Segment). In IEP-Ref, each of these modules is parameterized with a small fully convolutional network and belongs to one of the following 4 categories:

- **Preprocess:** This component maps the image to the feature tensor. Its output is the input to the Scene module. See Table 4 for the network architecture.
- **Unary:** This includes the Scene, Filter_X, Unique, Relate, Same_X modules. It transforms one feature tensor to another. See Table 5 for the network architecture.
- **Binary:** This includes the And and Or modules. It transforms two feature tensors to one. See Table 6 for the network architecture.
- **Postprocess:** This only includes the Segment module. It transforms the 128-channel feature tensor to a

Layer	Output size
Input image	$3 \times 320 \times 320$
ResNet101 [8] conv4_6	$1024 \times 20 \times 20$
Conv(3×3 , $1024 \rightarrow 128$)	$128 \times 20 \times 20$
ReLU	$128 \times 20 \times 20$
Conv(3×3 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
ReLU	$128 \times 20 \times 20$

Table 4: Network architecture for the **Preprocess** module.

Index	Layer	Output size
(1)	Previous module output	$128 \times 20 \times 20$
(2)	Conv(3×3 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
(3)	ReLU	$128 \times 20 \times 20$
(4)	Conv(3×3 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
(5)	Residual: Add (1) and (4)	$128 \times 20 \times 20$
(6)	ReLU	$128 \times 20 \times 20$

Table 5: Network architecture for the **Unary** modules.

Index	Layer	Output size
(1)	Previous module output	$128 \times 20 \times 20$
(2)	Previous module output	$128 \times 20 \times 20$
(3)	Concatenate (1) and (2)	$256 \times 20 \times 20$
(4)	Conv(1×1 , $256 \rightarrow 128$)	$128 \times 20 \times 20$
(5)	ReLU	$128 \times 20 \times 20$
(6)	Conv(3×3 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
(7)	ReLU	$128 \times 20 \times 20$
(8)	Conv(3×3 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
(9)	Residual: Add (5) and (8)	$128 \times 20 \times 20$
(10)	ReLU	$128 \times 20 \times 20$

Table 6: Network architecture for the **Binary** modules.

Layer	Output size
Previous module output	$128 \times 20 \times 20$
Unary module	$128 \times 20 \times 20$
Conv(1×1 , $128 \rightarrow 128$)	$128 \times 20 \times 20$
ReLU	$128 \times 20 \times 20$
Bilinear upsample	$128 \times 320 \times 320$
Conv(1×1 , $128 \rightarrow 128$)	$128 \times 320 \times 320$
ReLU	$128 \times 320 \times 320$
Conv(1×1 , $128 \rightarrow 32$)	$32 \times 320 \times 320$
ReLU	$32 \times 320 \times 320$
Conv(1×1 , $32 \rightarrow 4$)	$4 \times 320 \times 320$
ReLU	$4 \times 320 \times 320$
Conv(1×1 , $4 \rightarrow 1$)	$1 \times 320 \times 320$

Table 7: Network architecture for the Segment module.

1-channel segmentation mask. See Table 7 for the network architecture.

Network architectures for **Preprocess**, **Unary**, **Binary** are directly inherited from IEP [16].

B. More Model Analysis on CLEVR-Ref+

B.1. Number of Objects in a Scene

We suspect that the more objects in a scene, the harder for the model to carry out the referring reasoning steps. In Figure 9 we plot the performance of each model with respect to the number of objects in a scene. All models drop in performance when the number of objects increases, suggesting that the models tend to struggle when dealing with too many objects.

B.2. Schedule of Acquiring Reasoning Abilities

We are interested to see if throughout the training process, the network exhibit a schedule of acquiring various reasoning abilities (e.g. spatial reasoning, logic etc). From Figure 10, it seems that no such schedule was developed, and performance steadily increase across different referring

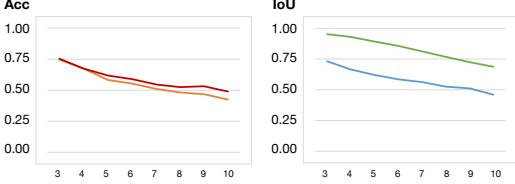


Figure 9: Effect of number of objects in a scene on referring detection or segmentation performance.

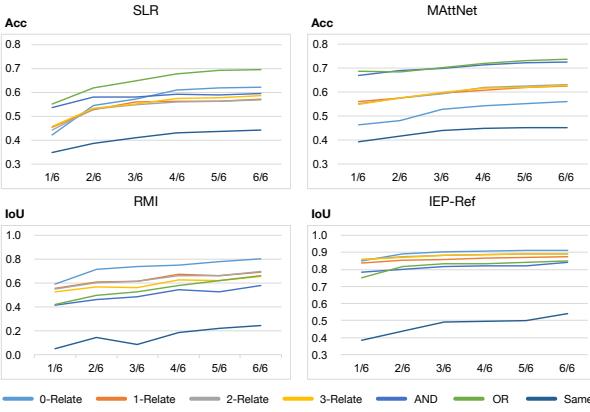


Figure 10: Performance across different referring expression categories throughout training. We inspect the performance every 1/6 of the entire training iterations.

expression categories. This may be due to the random sampling during training, instead of active learning (c.f. [28]).

B.3. Novel Compositions

To further test the models’ generalization ability, we also conducted experiments on the Compositional Generalization Test (CoGenT) data provided by CLEVR [15]. Here models are trained on objects with only a subset of all combinations, and then tested on both the same subset of combinations ($valA$) and another subset of combinations ($valB$). Results are summarized in Figure 11. We see a very small gap for detection models, suggesting that they have learned compositionality to generalize well. The gap for segmentation models, on the other hand, is larger.

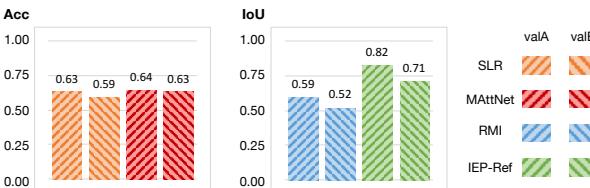
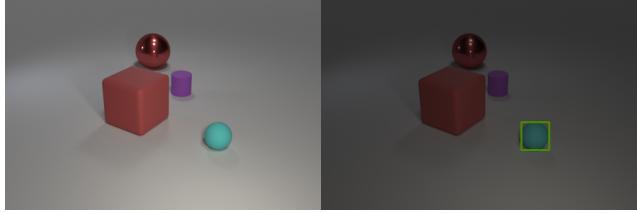


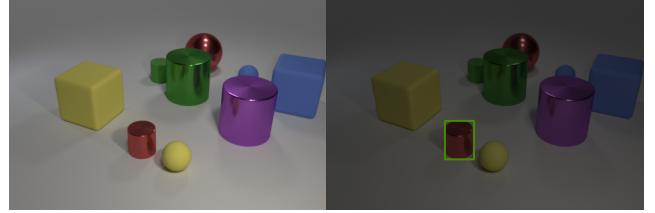
Figure 11: Different models’ performance on $valA$ and $valB$ of the CLEVR CoGenT data.

C. More Data Examples from CLEVR-Ref+

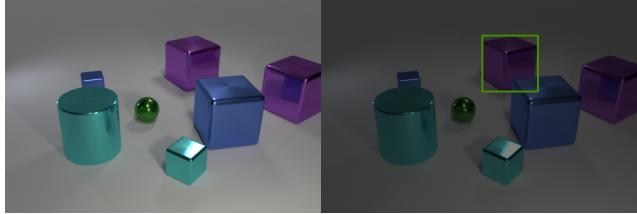
The remaining pages show random images, referring expressions, and the referring ground truth from our CLEVR-Ref+ dataset. In particular, we choose at least one example from each referring expression category (the 7 middle columns in Table 3 of the main paper). We show both detection ground truth (Figure 12) and segmentation ground truth (Figure 13).



(a) Look at matte thing that is on the left side of the red object that is behind the second one of the object(s) from right; The first one of the rubber thing(s) from front that are right of it



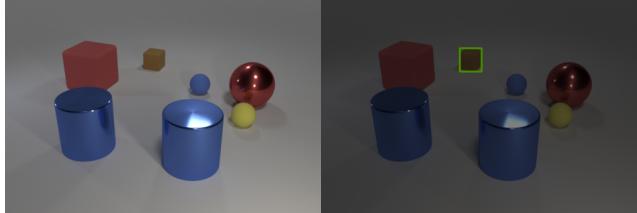
(b) The objects that are the seventh one of the thing(s) from right that are in front of the ninth one of the thing(s) from front or the second one of the thing(s) from front



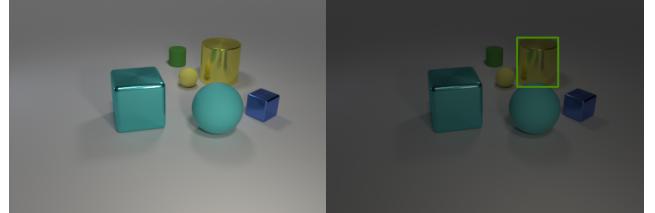
(c) The big metallic object(s) that are both to the left of the third one of the large thing(s) from left and on the right side of the first one of the object(s) from front



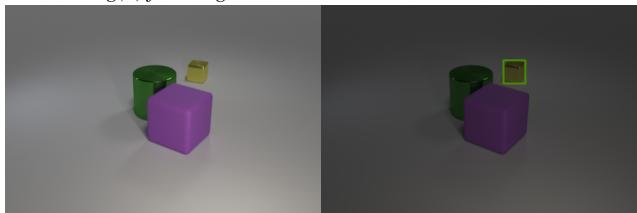
(d) The fully visible yellow ball(s)



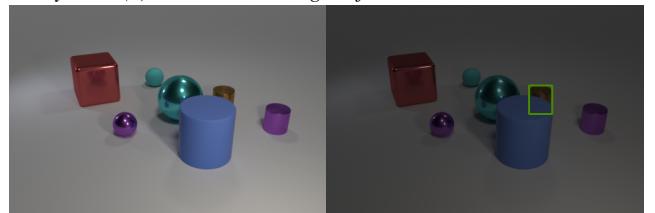
(e) Any other things that are the same shape as the fourth one of the rubber thing(s) from right



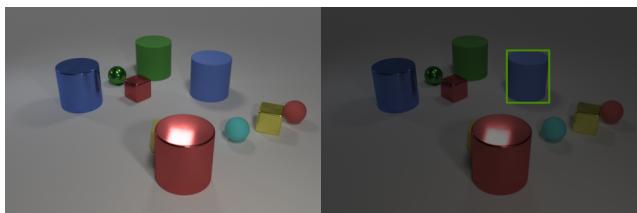
(f) Find object that is behind the fifth one of the object(s) from left; The cylinder(s) that are to the right of it



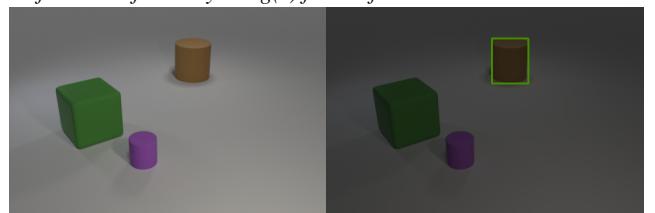
(g) Look at partially visible object(s); The second one of the thing(s) from left that are on the right side of it



(h) The second one of the shiny cylinder(s) from right that are to the right of the thing that is behind the thing that is on the left side of the first one of the tiny thing(s) from left

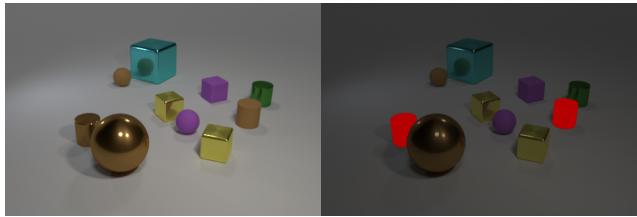


(i) The blue things that are either the fourth one of the thing(s) from right or the first one of the tiny ball(s) from front

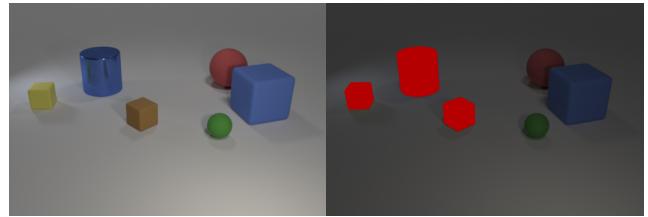


(j) The matte object(s) that are behind the second one of the cylinder(s) from right and on the right side of the first one of the object(s) from left

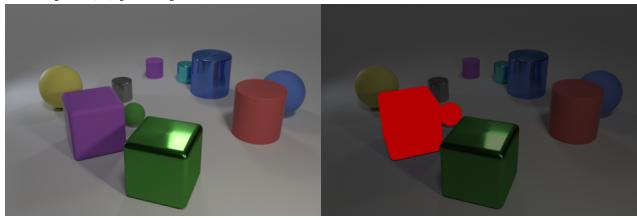
Figure 12: Referring object detection examples from CLEVR-Ref+.



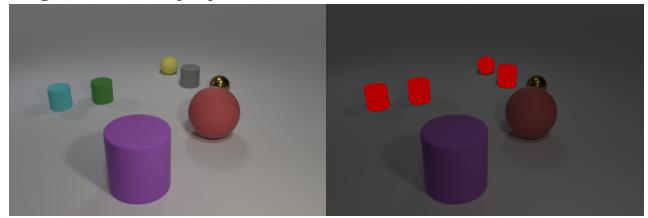
(a) Any other things that are the same shape as the seventh one of the object(s) from front



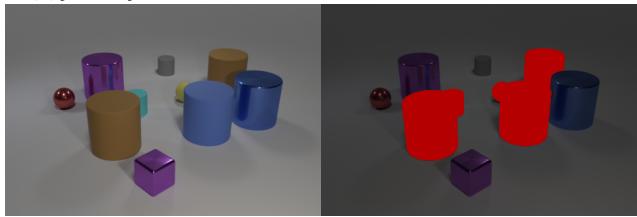
(b) Look at rubber ball that is to the left of the red ball(s); The thing(s) that are left of it



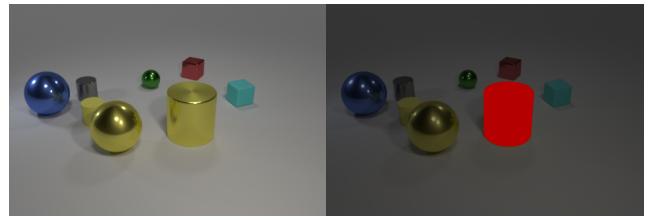
(c) The rubber object(s) that are to the right of the sixth one of the rubber thing(s) from right and to the left of the fifth one of the object(s) from left



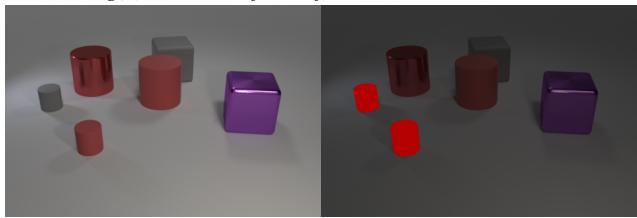
(d) The fully visible small thing(s)



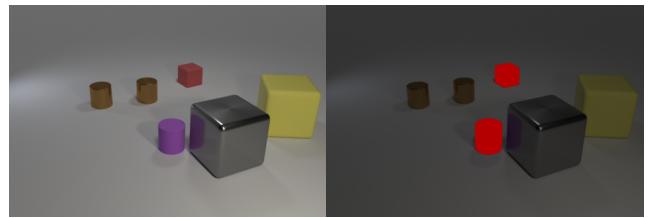
(e) Look at tiny rubber cylinder that is behind the tiny object that is on the right side of the seventh one of the cylinder(s) from front; The rubber thing(s) that are in front of it



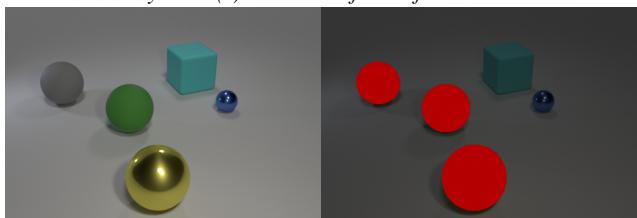
(f) The big things that are the sixth one of the object(s) from left or the seventh one of the object(s) from right



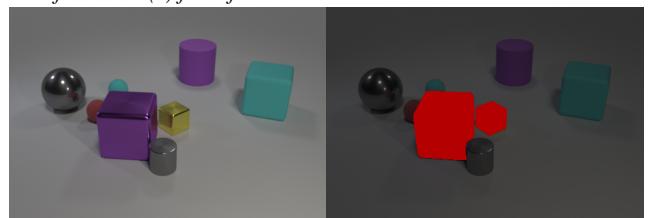
(g) Find the second one of the red rubber thing(s) from left; The fully visible rubber cylinder(s) that are in front of it



(h) Any other tiny object(s) made of the same material as the second one of the cube(s) from front



(i) Look at object that is to the right of the fourth one of the big object(s) from front; The ball(s) that are to the left of it



(j) The metallic object(s) that are behind the fourth one of the object(s) from right and in front of the fourth one of the thing(s) from front

Figure 13: Referring image segmentation examples from CLEVR-Ref+.