A. Appendix for "Meta Module Network for Compositional Visual Reasoning"

A.1. Visual Encoder and Multi-head Attention

The multi-head attention network is illustrated in Figure 7.



Figure 7. Illustration of the multi-head attention network used in the Meta Module.

A.2. Recipe Embedding

The recipe embedder is illustrated in Figure 8.



Figure 8. Illustration of the recipe embedder.

A.3. Implementation

The implementation of the proposed model is demonstrated in Figure 9. Our model can be efficiently implemented by adding masks on top of the Transformer model, guided by an additional supervision signal.



Figure 9. Illustration of the implementation of both the Transformer model (left) and our model (right).

A.4. Function Statistics

The function statistics is listed in Table 5.

Туре	Relate	Select	Filter	Choose	Verify	Query	Common	Differ	Bool	Exist	All
Funcs	5	1	8	12	5	6	2	6	2	1	48
Table 5. The statistics of different functions.											
Binary				Objects			String				
• •					1			-		_	
verify	choose	compare	ex1st	and	or fil	ter sele	ct relate	query[o	bject]	query[so	cene

Table 6. Error analysis on different types of functions. "Objects" functions only appear in the intermediate step, "String" function only appears in the final step, "Binary" functions can occur in both intermediate and final step

A.5. Inferential Chains

More inferential chains are visualized in Figure 10 and 11.



Figure 10. More examples on visualization of the inferential chains learned by our model.

A.6. Detailed Error Analysis

Furthermore, we conduct a detailed analysis of function-wise execution accuracy to understand the limitation of MMN. Results are shown in Table 6. Below are the observed main bottlenecks: (*i*) relation-type functions such as relate, relate_inv; and (*ii*) object/attribute recognition functions such as query_name, query_color. We hypothesize that this might be attributed to the quality of visual features from standard object detection models [1], which does not capture the relations between objects well. Besides, the object and attribute classification network are not fine-tuned on GQA. This suggests that scene graph modeling for visual scene understanding is critical to surpassing NSM [18] on performance.

A.7. Function Description

The detailed function descriptions for CLEVR and GQA are provided in Figure 12 and Figure 13, respectively.



What is the dark clothing? What is the man getting on? What do you think is the standing person near the man wearing ?

Figure 11. More examples on visualization of the inferential chains learned by our model.

Туре	Overrides	Arg0 ? Means Dependecy	arg1	Output
Relate	-	?	Position	Region
Filter	color	?	color	Region
	material	?	material	Region
	shape	?	shape	Region
	size	?	size	Region
Same	color	?	?	Region
	material	?	?	Region
	shape	?	?	Region
	size	?	?	Region
Query	color	?	color	Answer
	material	?	material	Answer
	shape	?	shape	Answer
	size	?	size	Answer
Logic	Union	?	?	Region
	Intersect	?	?	Region
Equal	color	?	?	Answer
	material	?	?	Answer
	shape	?	?	Answer
	size	?	?	Answer
Compare	greater_than	?	?	Answer
	less_than	?	?	Answer
Count	-	?	-	Answer
exist	-	?	-	Answer

Figure 12. The function definitions and their corresponding outputs on CLEVR.

Туре	Overrides	arg0 (? Means Dependency)	arg1	arg2	arg3	Output		
	Relate	?	Relation	-	-			
	Relate_with_name	?	Relation	Object	-			
Relationship	Relate_invese	?	Relation	-	-	Region		
	Relate_inverse_with_name	?	Relation	Object	-			
	Relate_with_same_attribute	?	Relation	Attribute	-			
Selection	Select	-	Object	-	-	Region		
	Filter_horizontal_position	?	H-Position	-	-			
	Filter_Vertical_position	?	V-Position	-	-	Region		
	Filter_with_color	?	Color	-	-			
	Filter_with_shape	?	Shape	-	-			
Filter	Filter_with_activity	?	Activity	-	-			
	Filter_with_material	?	Material	-	-			
	Filter with color noteq	?	Color	-	-			
	Filter with shape noteq	?	Shape	-	-			
	Choose name	?	Name1	Name2	-			
	 Choose scene	-	Scene1	Scene2	-			
	Choose color	?	Color1	Color2	-			
	Choose shape	?	Shape1	Shape2	-			
	Choose horizontal position	?	H-Position1	H-Position2	-	- - - Answer		
	Choose vertical position	2	V-Position1	V-Position2	-			
Choose	Choose relation name	?	Relation1	Relation2	Name			
	Choose relation inverse name	?	Relation1	Relation2	Name			
	Choose vounger	?	?	-	-			
	Choose older	2	?	-	-			
	Choose healthier	?	?	_	_			
	Choose less healthier	?	?	-	-			
	Verify color	2	Color	-	-	Answer		
	Verify shape	2	Shape	-	-			
Verify	Verify scene	-	Scene	_	_			
verny	Verify relation name	2	Belation	Name	_			
	Verify_relation_inv_name	?	Relation	Name	_			
		?	-	-	_	Answer		
		2						
		2						
Query		-						
		2						
	Query_nonzontal_position	?	_	-	_			
					_			
Common	Common material	[:::.:]				Answer		
	Different name	[:::]	-	-	-			
Different	Different_name	ر: ۲۰۰۱ ۲	-	-	-			
		2	2	-	-	Answer		
		י ני בי בי	ŗ	-	-	+		
Como		[[[[]]]]	-	-	-	Answer		
same	Same_name	r 2	۲ ۲	-	-			
	Same_color	· ·	1 2	-	-	<u> </u>		
And	And	? 	?	-	-	Answer		
Or	Or	?	?	-	-	Answer		
Exist	Exist	2	2	-	-	Answer		

Figure 13. The function definitions and their corresponding outputs on GQA.