

---

# Systematic Generalization: What Is Required and Can It Be Learned?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Numerous models for grounded language understanding have been recently pro-  
2 posed, including (i) generic models that can be easily adapted to any given task  
3 and (ii) intuitively appealing Neural Module Networks (Andreas et al., 2016a) that  
4 require background knowledge to be instantiated. We compare generic and modular  
5 models in how much they lend themselves to a particular form of systematic gener-  
6 alization. Our findings show that the generalization of modular models is much  
7 more systematic and that it is highly sensitive to the module layout, i.e. to how  
8 exactly the modules are connected. We furthermore investigate if modular models  
9 that generalize well could be made more end-to-end by learning their layout and  
10 parametrization. We show how end-to-end methods from prior work often learn  
11 spurious layouts and parametrizations that do not facilitate systematic generaliza-  
12 tion. Our results suggest that, in addition to modularity, systematic generalization  
13 in language understanding may require explicit regularizers or priors.

## 14 1 Introduction

15 In recent years, neural network based models have become the workhorse of natural language under-  
16 standing and generation showing state-of-the-art performance on numerous benchmarks, including  
17 Recognizing Textual Entailment (RTE) (Gong et al., 2017), Visual Question Answering (VQA)  
18 (Jiang et al., 2018), and Reading Comprehension (Wang et al., 2018). Despite these successes, a  
19 growing body of literature suggests that these approaches latch onto statistical regularities which are  
20 omnipresent in existing datasets (Agrawal et al., 2016; Gururangan et al., 2018; Jia & Liang, 2017)  
21 and do not generalize outside of the specific distributions they are trained on. These findings have  
22 recently been corroborated by Lake & Baroni (2018), who showed that seq2seq models (Sutskever  
23 et al., 2014; Bahdanau et al., 2015) show little systematicity (Fodor & Pylyshyn, 1988) in how they  
24 generalize, i.e. they fail to learn general rules on how to compose words.

25 Introduced by Andreas et al. (2016b), Neural Module Networks (NMNs) approach aims to improve  
26 the generalization capabilities of neural models by adding *modularity* and *structure* to their design  
27 to make them resemble the kind of rules they are supposed to learn. The NMN approach, while  
28 intuitively appealing, has seen limited adoption because of the large amount of domain knowledge it  
29 requires to decide (Andreas et al., 2016a) or predict (Johnson et al., 2017; Hu et al., 2017) how the  
30 modules should be created (parametrization) and how they should be connected (layout) based on a  
31 query. Besides, their performance has often been matched by more generic neural models, such as  
32 FiLM (Perez et al., 2017), Relations Networks (Santoro et al., 2017), and CAN (Hudson & Manning,  
33 2018), and their generalization, to the best of our knowledge, has not been a subject of a focused  
34 study. In this work we investigate the impact of explicit modularity and structure on systematic  
35 generalization by studying the generalization of NMNs and contrasting it to those of generic models.  
36 We choose to focus on the following basic generalization requirement: *a good model should be able*

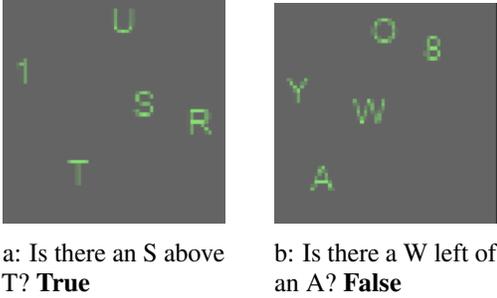


Figure 1: SGOOP tests model’s ability to reason about all object pairs after being trained on a small subset of them. We show a positive (left) and negative (right) example from SGOOP.

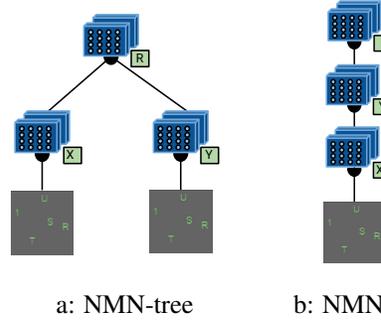


Figure 2: The different NMN layouts for the question “is there a S above T”.

37 to reason about all possible object combinations despite being trained on a small subset of them. We  
 38 instantiate this requirement in form of a simple yes-no Visual Question Answering (VQA) dataset.

39 Our first finding is that NMNs do generalize better than other neural models when an appropriate  
 40 choice of layout and parametrization is made. We furthermore experiment with existing methods for  
 41 making NMNs more end-to-end by inducing the module layout (Johnson et al., 2017) or learning  
 42 module parametrization through soft-attention over the question (Hu et al., 2017). We show how  
 43 such end-to-end approaches often fail to find the right structural settings and instead prefer a wrong  
 44 chain layout or spurious parametrization and do not generalize better than the generic models. We  
 45 believe that our findings challenge the intuition of researchers in the field and provide a foundation  
 46 for improving systematic generalization of neural approaches to language understanding.

## 47 2 Setup

48 **Dataset:** SGOOP (Spatial Queries Over Object Pairs, Figure 1) is a minimalistic VQA task designed  
 49 to test a particular type of generalization: the ability to disentangle the meaning of relation words  
 50 and object words and then compose these meanings in novel contexts to perform basic relational  
 51 reasoning in a consistent way. Concretely, SGOOP requires answering a yes-no question  $q = X R Y$   
 52 about whether objects  $X$  and  $Y$  are in a spatial relation  $R$ , given a  $64 \times 64$  RGB image  $x$ .  $x$  contains  
 53 5 randomly chosen and randomly positioned objects. There are 36 objects: letters A-Z and digits 0-9,  
 54 and 4 relations: LEFT\_OF, RIGHT\_OF, ABOVE, and BELOW. Our goal is to discover which models  
 55 can correctly answer questions about all  $36 \cdot 36$  possible object pairs in the SGOOP dataset after  
 56 having been trained on only a subset. Therefore, we train on  $36 \cdot 4 \cdot k$  unique questions, where for  
 57 every left-hand-side (LHS) object  $X$ , we randomly sample  $k$  different right-hand-side objects (RHS),  
 58 and test on the remaining  $36 \cdot 4 \cdot (36 - k)$  questions. We refer to  $k$  as the *#rhs/lhs parameter* of the  
 59 dataset. To exclude a possible compounding factor of overfitting on training images, all our training  
 60 sets contain 1 million examples obtained by sampling multiple images per question.

61 **Models:** We experiment with models from 2 broad categories. *Generic* models such as FiLM (Perez  
 62 et al., 2017), Relation Networks (RelNet, Santoro et al. (2017)) and CAN Hudson & Manning  
 63 (2018), and *modular* and *structured* Neural Module Networks (NMN). NMNs (Andreas et al., 2016b)  
 64 construct question-specific networks by composing together trainable neural modules. To answer  
 65 a question with an NMN, a computation graph is constructed by making 2 decisions: *layout* -  
 66 the number of modules, their types and how they are connected, and *parametrization* - how these  
 67 modules are parametrized based on the question. For our study we adapt the N2NMN (Hu et al., 2017)  
 68 paradigm from this family, which predicts the layout with a seq2seq model (Sutskever et al., 2014)  
 69 and computes the parametrization of the modules using a soft attention mechanism. Since all the  
 70 questions in SGOOP have the same structure, we can get away with a single trainable layout variable  
 71 and separate trainable attention variables per each module. We also experiment with hard-coded  
 72 layout and parametrization setting, in the spirit of original NMN (Andreas & Klein, 2015).

73 Formally, our NMN is constructed by repeatedly applying a *generic neural module*  $f(\theta, \gamma, h_l, h_r)$ ,  
 74 which takes as inputs the shared parameters  $\theta$ , the question-specific parametrization  $\gamma$  and the left-

Table 1: Comparing the performance of generic models to the structured NMN-Tree model on the hardest version of our dataset (lower #rhs/lhs is more difficult).

model	train. acc (%)	test acc. (%)
Conv+LSTM	97.9	64.4 ± 1.8
RelNet	95.6	63.1 ± 1.0
FiLM	100	66.6 ± 2.5
MAC	99.5	72.6 ± 3.4
NMN-Tree (Residual)	100	100.0 ± 0.0
NMN-Tree (Find)	100.0	99.7 ± 0.3
NMN-Chain (Find)	99.2	51.4 ± 2.8
NMN-Chain-XYR (Residual)	100	51.6 ± 1.6
NMN-Chain-XYR (Residual)	99.7	54.1 ± 1.7
NMN-Chain-RXY (Residual)	98.7	50.5 ± 0.9

75 hand side and right-hand side inputs  $h_l$  and  $h_r$ .  $M$  such modules are connected and conditioned on a  
 76 question  $q = (q_1, q_2, q_3)$  as follows:

$$\gamma_k = \sum_{i=1}^s \alpha^{k,i} e(q_i) \quad (1)$$

$$h_k = f(\theta, \gamma_k, \sum_{j=-1}^{k-1} \tau_0^{k,j} h_j, \sum_{j=-1}^{k-1} \tau_1^{k,j} h_j) \quad (2)$$

77 In the equations above,  $h_{-1} = 0$  is the zero tensor,  $h_0 = h_x$  are the image features outputted by a  
 78 CNN network referred to as the *stem*, and  $e$  is the embedding table for the questions words. We refer  
 79 to  $A = (\alpha^{k,i})$  and  $T = (\tau_0^{k,i}, \tau_1^{k,i})$  as the *parametrization attention matrix* and the *layout tensor*  
 80 respectively. The output of the final module,  $h_M$  is fed into a fully connected classifier network to  
 81 make predictions. We perform our experiments with the Find module from Hu et al. (2017) and the  
 82 Residual module from Johnson et al. (2017) with minor modifications (Appendix-A.3 for details).

83 Based on the generic NMN model described above, we experiment with several specific architectures.  
 84 Each of the models uses  $M = 3$  modules, which are connected and parametrized differently. In  
 85 **NMN-Chain** the modules form a sequential chain as shown in Figure 1b. Modules 1, 2 and 3 are  
 86 parametrized based on the first object word, second object word and the relation word respectively,  
 87 which is achieved by setting the attention  $\alpha_1, \alpha_2, \alpha_3$  to the corresponding one-hot vectors. **NMN-**  
 88 **Tree** has similar hard-coded attention vectors, with a tree-like connectivity between the modules.  
 89 **1a**. In the **Stochastic N2NMN**, similar to the N2NMN (Hu et al., 2017), the layout  $T$  is treated as a  
 90 stochastic latent variable that takes two values:  $T_{tree}$  as in NMN-Tree, and  $T_{chain}$  as in NMN-Chain.  
 91 The output probabilities are computed by marginalizing out  $T$ , i.e. probability of label “yes” is  
 92 computed as  $p(\text{yes}|x, q) = \sum_{T \in \{T_{tree}, T_{chain}\}} p(\text{yes}|T, x, q)p(T)$ . **Attention N2NMN** also from  
 93 (Hu et al., 2017), is structured just like NMN-Tree but has  $\alpha^k$  computed as  $\text{softmax}(\tilde{\alpha}^k)$ , where  $\tilde{\alpha}^k$   
 94 is a trainable vector. We use Attention N2NMN only with the Find module, which was designed by  
 95 (Hu et al., 2017) specifically parametrized with the help of soft attention.

### 96 3 Experiments

97 **Which Models Generalize Better:** We report the performance for all models on the hardest version  
 98 of our dataset (#rhs/lhs = 1) in Table 1. These results show that generic models do not generalize  
 99 well, while the NMN-Tree model does. To understand better the cause of NMN-Tree’s advantage  
 100 we compare the performance of the NMN-Tree and NMN-Chain models. The results show that for  
 101 both Find and Residual architectures, using a tree layout is crucial for generalization as NMN-Chain  
 102 performs barely above random chance.

103 **Can the Right Kind of NMN Be Induced:** The generalization of NMN-Tree model, while impres-  
 104 sive, is somewhat unsurprising because both the *layout* and *parametrization* of this model encode a  
 105 significant amount of prior knowledge about the task. We therefore investigate whether the amount  
 106 of such prior knowledge can be reduced by fixing one of the structural aspects and inducing the other.  
 107 For inducing a layout, we use the Stochastic N2NMN model. We experiment with both Find and

Table 2: Layout induction results for Stochastic N2NMNs using Residual modules and Find modules. For each setting of  $p_0(tree)$ , we report results on 1 rhs/lhs and 18 rhs/lhs datasets.

(a) Residual modules				(b) Find modules			
#rhs/lhs	$p_0(tree)$	test acc. (%)	$p_{50K}(tree)$	#rhs/lhs	$p_0(tree)$	test acc. (%)	$p_{50K}(tree)$
1	0.1	$52.7 \pm 2.2$	0.003	1	0.1	$51.2 \pm 2.9$	0.00
	0.5	$57.0 \pm 4.4$	0.026		0.5	$93.2 \pm 7.1$	0.999
	0.9	$99.9 \pm 0.1$	0.997		0.9	$95.9 \pm 1.6$	0.999
18	0.1	$100.0 \pm 0.0$	0.999	18	0.1	$78.6 \pm 20.7$	0.2
	0.5	$97.7 \pm 5.1$	0.999		0.5	$91.6 \pm 6.5$	0.999
	0.9	$99.1 \pm 2.3$	0.999		0.9	$97.3 \pm 3.4$	0.999

108 Residual modules and report results with diverse initial conditions,  $p_0(tree) = 0.1, 0.5, 0.9$ , where  
 109  $p_0(tree)$  is the initial value of  $p(T_{tree})$ . The results obtained on the #rhs/lhs=1 dataset (Table 2) show  
 110 that the correct layout was not induced for  $p_0(tree) = 0.1$  and  $p_0(tree) = 0.5$  with the Residual  
 111 module and for  $p_0(tree) = 0.1$  with the Find module. We also run similar experiments on an easy-to-  
 112 generalize #rhs/lhs=18 version. Here, the NMN with Residual module preferred the tree layout for  
 113 all initializations. It is notable, however, that in the setting with #rhs/lhs=1 where the correct choice  
 114 of the layout is the *only* way to generalize, only a very lucky initialization  $p_0(tree) = 0.9$  resulted in  
 115 successful layout induction for the Residual module.

116 For parameterization induction, we experiment with the Attention N2NMN model on #rhs/lhs=1  
 117 and #rhs/lhs=18. The model often did not find the attention settings that lead to generalization on  
 118 the challenging #rhs/lhs=1 split (83.8% test accuracy). That should be contrasted with the close-to-  
 119 perfect 99.2% accuracy of the model that was trained on #rhs/lhs=18 version of the task, suggesting  
 120 that the parametrization induction did not work due to the difficulty of our #rhs/lhs=1 split. To analyze  
 121 the learnt attention weights, we compute a sharpness ratio  $\rho = \max(\alpha^{k,X}, \alpha^{k,Y}) / \min(\alpha^{k,X}, \alpha^{k,Y})$   
 122 for modules  $k = 1$  and  $k = 2$  for each of the trained modules. We find that the learnt attention weights  
 123 on #rhs/lhs=1 are generally blurry with  $\rho < 2$  for 40% of the modules (details in Appendix-A.4).

## 124 4 Related Work

125 Using synthetic VQA datasets to study grounded language understanding is a recent trend started by  
 126 CLEVR (Johnson et al., 2016), and recently the ShapeWorld dataset (Kuhnle & Copestake, 2017),  
 127 that involves a number of VQA generalization tests. Closely related to our work is the recent study  
 128 on generalization to long-tail questions about rare objects by Bingham et al. (2018). They do not,  
 129 however, consider as many models as we do and do not study whether the best-performing models  
 130 can be made end-to-end. Andreas et al. (2016a) introduced NMNs as a modular, structured VQA  
 131 model where a fixed number of hand-crafted neural modules are chosen and composed together in a  
 132 layout determined by the dependency parse of the question. Hu et al. (2017) and Johnson et al. (2017)  
 133 followed up with end-to-end NMNs, removing the non-differentiable parser. Recent concurrent work  
 134 by (Hu et al., 2018) attempts to remove the need for hard stochastic layout decisions.

## 135 5 Conclusion and Discussion

136 We have conducted a rigorous investigation of an important form of systematic generalization required  
 137 for grounded language understanding: the ability to reason about all possible pairs of objects despite  
 138 being trained on a small subset. The intuitive appeal of modularity and structure in designing  
 139 neural architectures for language is now supported by our results. Our other key finding is that  
 140 coming up with an end-to-end and/or soft version of modular models may be not sufficient for  
 141 strong generalization, because in the very setting where strong generalization is required, end-to-end  
 142 methods may find a different, less compositional solution (e.g. a chain layout or blurred attention).  
 143 This conclusion is relevant in the view of recent work done in the direction of making Neural Module  
 144 Networks more end-to-end (Suarez et al., 2018; Hu et al., 2018; Hudson & Manning, 2018). We hope  
 145 that our findings will inform researchers working on language understanding and provide them with  
 146 a useful intuition about what facilitates strong generalization and what is likely to inhibit it.

147 **References**

- 148 Aishwarya Agrawal, Dhruv Batra, and Devi Parikh. Analyzing the Behavior of Visual Question  
149 Answering Models. *arXiv:1606.07356 [cs]*, June 2016. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1606.07356)  
150 [1606.07356](http://arxiv.org/abs/1606.07356). arXiv: 1606.07356.
- 151 Jacob Andreas and Dan Klein. Alignment-Based Compositional Semantics for Instruction Following.  
152 In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*,  
153 2015. URL <https://arxiv.org/abs/1508.06491>.
- 154 Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. *2016*  
155 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 39–48, 2016a.
- 156 Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural Module Networks. In  
157 *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
158 2016b. URL <http://arxiv.org/abs/1511.02799>.
- 159 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly  
160 Learning to Align and Translate. In *Proceedings of the ICLR 2015*, 2015.
- 161 Eli Bingham, Piero Molino, Paul Szerlip, Obermeyer Fritz, and Goodman Noah. Chacterizing how  
162 Visual Question Answering scales with the world. In *NIPS Workshop on Visually-Grounded*  
163 *Interaction and Language*, 2018.
- 164 Jerry A. Fodor and Zenon W. Pylyshyn. Connectionism and cognitive architecture: A critical analysis.  
165 *Cognition*, 28(1):3–71, 1988.
- 166 Yichen Gong, Heng Luo, and Jian Zhang. Natural Language Inference over Interaction Space.  
167 *arXiv:1709.04348 [cs]*, September 2017. URL <http://arxiv.org/abs/1709.04348>.  
168 arXiv: 1709.04348.
- 169 Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and  
170 Noah A. Smith. Annotation Artifacts in Natural Language Inference Data. *arXiv:1803.02324 [cs]*,  
171 March 2018. URL <http://arxiv.org/abs/1803.02324>. arXiv: 1803.02324.
- 172 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
173 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
174 pp. 770–778, 2016.
- 175 Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to  
176 Reason: End-to-End Module Networks for Visual Question Answering. *arXiv:1704.05526 [cs]*,  
177 April 2017. URL <http://arxiv.org/abs/1704.05526>. arXiv: 1704.05526.
- 178 Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko. Explainable neural computation via  
179 stack neural module networks. In *Proceedings of the European Conference on Computer Vision*  
180 *(ECCV)*, pp. 53–69, 2018.
- 181 Drew A. Hudson and Christopher D. Manning. Compositional Attention Networks for Machine  
182 Reasoning. In *ICLR 2018*, February 2018.
- 183 Robin Jia and Percy Liang. Adversarial Examples for Evaluating Reading Comprehension Systems.  
184 *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*,  
185 pp. 2021–2031, 2017. doi: 10.18653/v1/D17-1215. URL [https://aclanthology.coli.](https://aclanthology.coli.uni-saarland.de/papers/D17-1215/d17-1215)  
186 [uni-saarland.de/papers/D17-1215/d17-1215](https://aclanthology.coli.uni-saarland.de/papers/D17-1215/d17-1215).
- 187 Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh.  
188 Pythia v0.1: The winning entry to the vqa challenge 2018. [https://github.com/](https://github.com/facebookresearch/pythia)  
189 [facebookresearch/pythia](https://github.com/facebookresearch/pythia), 2018.
- 190 Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and  
191 Ross Girshick. CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual  
192 Reasoning. *arXiv:1612.06890 [cs]*, December 2016. URL [http://arxiv.org/abs/1612.](http://arxiv.org/abs/1612.06890)  
193 [06890](http://arxiv.org/abs/1612.06890). arXiv: 1612.06890.

- 194 Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Judy Hoffman, Li Fei-Fei, C. Lawrence  
195 Zitnick, and Ross Girshick. Inferring and Executing Programs for Visual Reasoning. In *ICCV*,  
196 2017. URL <http://arxiv.org/abs/1705.03633>.
- 197 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980*  
198 *[cs]*, December 2014. URL <http://arxiv.org/abs/1412.6980>. arXiv: 1412.6980.
- 199 Alexander Kuhnle and Ann Copestake. ShapeWorld - A new test methodology for multimodal  
200 language understanding. *arXiv:1704.04517 [cs]*, April 2017. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1704.04517)  
201 [1704.04517](http://arxiv.org/abs/1704.04517). arXiv: 1704.04517.
- 202 Brenden M. Lake and Marco Baroni. Generalization without systematicity: On the compositional  
203 skills of sequence-to-sequence recurrent networks. In *ICML*, 2018.
- 204 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual  
205 Reasoning with a General Conditioning Layer. In *In Proceedings of the AAAI Conference on*  
206 *Artificial Intelligence*, 2017. URL <http://arxiv.org/abs/1709.07871>.
- 207 Adam Santoro, David Raposo, David G. T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter  
208 Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning.  
209 *arXiv:1706.01427 [cs]*, June 2017. URL <http://arxiv.org/abs/1706.01427>. arXiv:  
210 1706.01427.
- 211 Joseph Suarez, Justin Johnson, and Fei-Fei Li. DDRprog: A CLEVR Differentiable Dynamic  
212 Reasoning Programmer. *arXiv:1803.11361 [cs]*, March 2018. URL [http://arxiv.org/](http://arxiv.org/abs/1803.11361)  
213 [abs/1803.11361](http://arxiv.org/abs/1803.11361). arXiv: 1803.11361.
- 214 Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to Sequence Learning with Neural Networks.  
215 In *Advances in Neural Information Processing Systems 27*, pp. 3104–3112, 2014.
- 216 Wei Wang, Ming Yan, and Chen Wu. Multi-Granularity Hierarchical Attention Fusion Networks for  
217 Reading Comprehension and Question Answering. In *Proceedings of the 56th Annual Meeting*  
218 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1705–1714,  
219 Melbourne, Australia, 2018. Association for Computational Linguistics. URL [http://aclweb.](http://aclweb.org/anthology/P18-1158)  
220 [org/anthology/P18-1158](http://aclweb.org/anthology/P18-1158).

221 **A Appendix**

222 **A.1 Additional Details**

223 **Dataset:** To make negative examples in SGOOP challenging, we ensure that both X and Y of a  
 224 question are always present in the associated image and that there are always distractor objects  
 225  $Y' \neq Y$  and  $X' \neq X$  such that  $XY'$  and  $X'Y$  are both true for the image. These extra  
 226 precautions guarantee that answering a question requires the model to locate all possible X and Y  
 227 then check if any pair of them are in the relation R.

228 **Hyperparameters:** All models share the same stem architecture which is a CNN based architecture  
 229 of 6 layers. Each layer is a Conv  $\rightarrow$  BatchNorm  $\rightarrow$  ReLU with a MaxPool after layers 1 and 3.  
 230 The input to the stem is a  $64 \times 64 \times 3$  image, and the feature dimension used throughout the stem is  
 231 64. All models are optimized using Adam [Kingma & Ba \(2014\)](#) with a learning rate of  $3e-4$ , and  
 232 with minibatches of size 128.

233 **A.2 Generic Models**

234 We consider four generic models in this paper: CNN+LSTM, FiLM, Relation Networks (RelNet),  
 235 and Compositional Attention Networks (CAN). For CNN+LSTM, FiLM, and RelNet models, the  
 236 question  $q$  is first encoded into a fixed-size representation  $h_q$  using a unidirectional LSTM network.

237 **CNN+LSTM** flattens the 3D tensor  $h_x$  to a vector and concatenates it with  $h_q$  to produce  $h_{qx}$ .

$$h_{qx} = [vec(h_x); h_q] \quad (3)$$

238 **RelNet** uses a network  $g$  which is applied to all pairs of feature columns of  $h_x$  concatenated with the  
 239 question representation  $h_q$ , all of which is then pooled to obtain  $h_{qx}$ :

$$h_{qx} = \sum_{i,j} g(h_x(i), h_x(j), h_q) \quad (4)$$

240 where  $h_x(i)$  is the  $i$ -th feature column of  $h_x$ .

241 **FiLM** networks use  $N$  convolutional FiLM blocks applied to  $h_x$ . A FiLM block is a residual block  
 242 ([He et al., 2016](#)) in which a feature-wise affine transformation (FiLM layer) is inserted after the 2nd  
 243 convolutional layer. The FiLM layer is conditioned on the question at hand via prediction of the  
 244 scaling and shifting parameters  $\gamma_n$  and  $\beta_n$ :

$$[\gamma_n; \beta_n] = W_q^n h_q + b_q^n \quad (5)$$

$$\tilde{h}_{qx}^n = BN(W_2^n * ReLU(W_1^n * h_{qx}^{n-1} + b_n)) \quad (6)$$

$$h_{qx}^n = h_{qx}^{n-1} + ReLU(\gamma_n \odot \tilde{h}_{qx}^n \oplus \beta_n) \quad (7)$$

245 where  $BN$  stands for batch normalization,  $*$  stands for convolution and  $\odot$  stands for element-wise  
 246 multiplications.  $h_{qx}^n$  is the output of the  $n$ -th FiLM block and  $h_{qx}^0 = h_x$ . The output of the last FiLM  
 247 block  $h_{qx}^N$  undergoes an extra  $1 \times 1$  convolution and max-pooling to produce  $h_{qx}$ .

248 **CAN** networks of [Hudson & Manning \(2018\)](#) produces  $h_{qx}$  by repeatedly applying a Memory-  
 249 Attention-Control (MAC) cell that is conditioned on the question through an attention mechanism.  
 250 The CAN model is quite complex and we refer the reader to the original paper for details.

251 **A.3 NMN Modules**

252 As mentioned in the text, our experiments are performed with the Find module from [Hu et al.](#)  
 253 [\(2017\)](#) and the Residual module from [Johnson et al. \(2017\)](#) with very minor modifications - we use  
 254 64 dimensional CNNs in our Residual blocks since our dataset consists of  $64 \times 64$  images. The  
 255 equations for the Residual module are as follows:

$$\theta = \emptyset, \quad (8)$$

$$\gamma = [W_1; b_1; W_2; b_2; W_3; b_3], \quad (9)$$

$$\tilde{h} = ReLU(W_3 * [h_l; h_r] + b_3), \quad (10)$$

$$f_{Residual}(\gamma, h_l, h_r) = ReLU(\tilde{h} + W_1 * ReLU(W_2 * \tilde{h} + b_2)) + b_1, \quad (11)$$

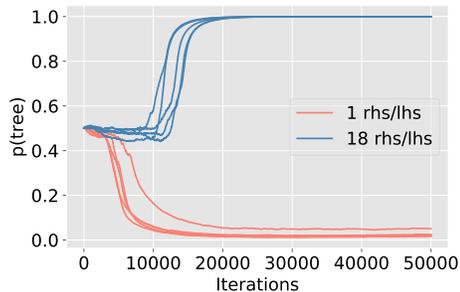


Figure 3: Learning dynamics of layout induction on 1 rhs/lhs and 18 rhs/lhs datasets using the Residual module with  $p_0(tree) = 0.5$ . All 5 runs of the model do not learn to use the tree layout for 1 rhs/lhs, the very setting where the tree layout is necessary for generalization.

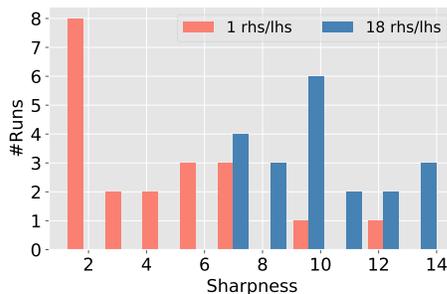


Figure 4: Histogram of sharpness ( $\rho$ ) values for attention weights induced on the 1 rhs/lhs and 18 rhs/lhs datasets. We can observe that the attention is much sharper for 18 rhs/lhs.

256 and for Find module as follows:

$$\theta = [W_1; b_1; W_2; b_2], \quad (12)$$

$$f_{Find}(\gamma, h_l, h_r) = ReLU(W_1 * \gamma \odot ReLU(W_2 * [h_l; h_r] + b_2) + b_1). \quad (13)$$

257 In formulas above  $W_1, W_2, W_3$  are convolution weights, and  $b_1, b_2, b_3$  are biases. The main difference  
 258 between Residual and Find is that in Residual all parameters depend on the questions words, where  
 259 as in Find convolutional weights are the same for all questions, and only the element-wise multipliers  
 260  $\gamma$  vary based on the question.

#### 261 A.4 Additional Results

262 **Structure Induction:** We visualize the progress of structure induction for the Residual module with  
 263  $p_0(tree) = 0.5$  in Figure 3. The figure shows  $p(tree)$  saturates to 0.0 or 1.0 eventually in #rhs/lhs=1  
 264 and #rhs/lhs=18 settings respectively.

265 **Parametrization Induction:** Figure 5 shows how attention weights evolve for an Attention N2NMN  
 266 model in the same context. It is notable that unlike in the gold-standard NMN-Tree model, the  
 267 relation word is mixed with the object words for modules 1 and 2. We also noticed that the model did  
 268 not learn to focus modules 1 and 2 on different words in the #rhs/lhs=1 setting (Figure 5b) as sharply  
 269 as it did in #rhs/lhs=18 (Figure-5a). To substantiate this observation with quantitative results, we  
 270 compute a sharpness ratio  $\rho = \max(\alpha^{k,X}, \alpha^{k,Y}) / \min(\alpha^{k,X}, \alpha^{k,Y})$  for modules  $k = 1$  and  $k = 2$   
 271 for each of the 20 modules that we have trained. One can observe from the histogram in Figure 4 that  
 272 attention weights learnt on #rhs/lhs=1 are generally blurry, with  $\rho$  being less than 2 for 8 modules out  
 273 of 20.

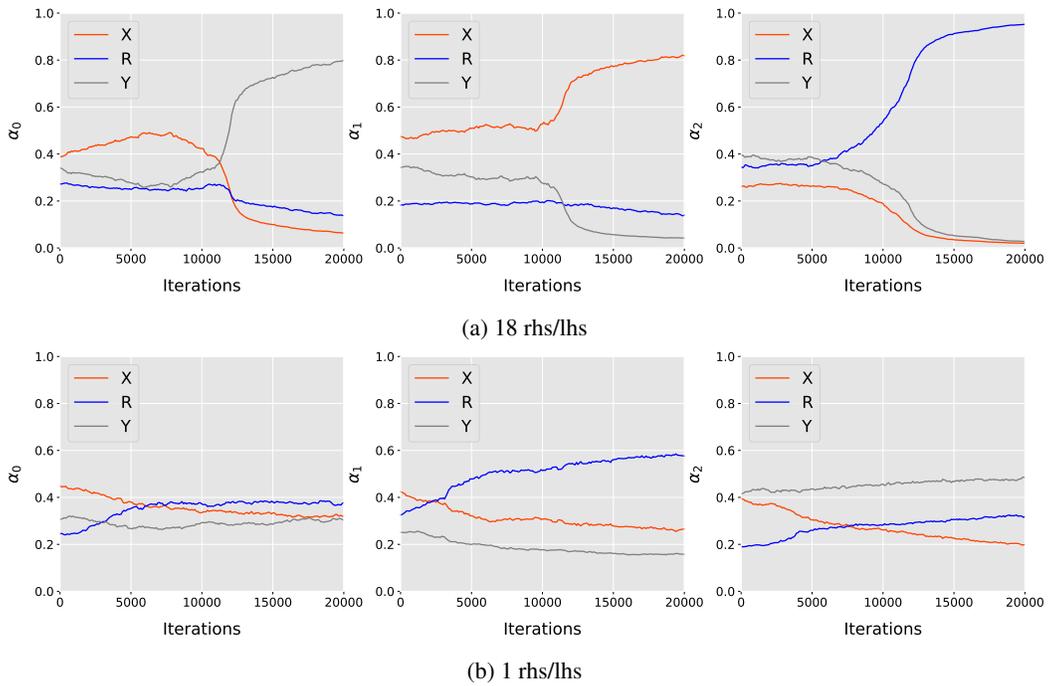


Figure 5: All three modules' attention weights for parametrization of the three question words for (a) 18 rhs/lhs and (b) 1 rs/lhs version of SQQOP. The model learns to disentangle between X and Y much better with more rhs/lhs.