

Dissecting Compositional Generalization: Correlation Analysis on Generalization Capacities of Different Compositional Problems

Yoobin Cheong, Yeong Koh and Yoon Tae Park

Center for Data Science

New York University

60 5th Avenue, New York, NY

{yc5206, yk2678, yp2201}@nyu.edu

1 Introduction

There exist many benchmark tests that claim to evaluate the abstract capacity of compositional generalization. However, whether these tests that fall under this abstraction target the same underlying capacity is an open question. (Weißenhorn et al., 2022; Furrer et al., 2020; Li et al., 2019; Liu et al., 2021; Hupkes et al., 2019)

As a means to investigate this idea, we will perform a correlation analysis on the performance of various neural network models trained on multiple tasks that have been proposed to measure compositional generalization. As a potential outcome of our project, grouping of tasks that are identified as being highly correlated will provide insights towards a better characterization of the abstraction compositional generalization, and furthermore be useful for task selection in multitask learning or transfer learning.

2 Experimental Design

Our team will start by finding relevant datasets and models to test. We will test some Sequence-to-Sequence models on widely used benchmark datasets such as SCAN (Lake and Baroni, 2018), CFQ (Keysers et al., 2019), and COGS (Kim and Linzen, 2020). After achieving good performances on given models, we will create a training evaluation pipeline, so that we can further test different models with different datasets or tasks. In this experiment, each task means a different split or generalization subcase of different datasets.

Then, we will compare the performances of the models through correlation analysis. For example, a high correlation between two tasks would be an indication that a network that performs well on one task performs well on another, suggesting that the two tasks recruit similar capacities. On the other hand, a low or no correlation would mean that they

may require distinct solutions as a solution for one task does not generalize well to another. At this stage, we will be working on the Greene cluster, to get time efficiency on multiple tests that would be conducted for next step and further developments.

3 Future Extension

This research can be further developed by applying the same experiments to other relevant datasets and models and analyzing the results. If the result shows some highly correlation between two tasks, we will also examine whether this data can be useful for task selection in multitask learning or transfer learning. If the result shows low or no correlation between two tasks, we will conclude that there are no straightforward correlations among different tasks based on our research and further investigation is needed.

4 Collaboration Statement

Our team will work altogether for all tasks in general, but will focus on different parts: (1) Yoobin Cheong will identify a set of compositional generalization datasets and models to test, and organize the different datasets into a single unified format, (2) Yoon Tae Park will build a codebase for the training evaluation pipeline, and (3) Yeong Koh will perform correlation analysis with the evaluation results and document any interesting findings. If correlation analysis shows some coherent, interpretable grouping of tasks, we will also test if this information can be helpful for task selection in multitask learning or transfer learning.

Acknowledgments

This research is supervised and mentored by Najeon Kim, a faculty fellow at CDS, as a part of her ongoing research project.

5 Experiment Progress

In response to the feedback we got on our proposal, we have clearly updated the structure of our experiment as follows:

- **What are the models you’re planning to train?** We are training on a T5-small Not-Pretrained model to compare Exact Match scores on different dataset/tasks. As there are many datasets within each task, our focus is to first experiment with the different datasets rather than various models. However, we plan to further explore other models, as time allows.
- **What architectures are you going to train?** We’ve created a pipeline that loads the data, encodes with a defined tokenizer and trains (finetunes) on the train set. We then evaluated the test set using Exact Match scores.
- **How are you going to get a large enough number of models to be able to compute correlations?** As mentioned, our research experiment is focused on different datasets/tasks. Therefore, as a baseline experiment, we’ve selected the T5 Not-Pretrained model and SCAN-simple/length/addprim turn left dataset/tasks to compute correlations. To enhance consistency, we are using 10 random seeds for every experiment, but will expand to more as time allows.
- **Are those going to be different architectures, different sets of hyperparameters, different random seeds?** We are setting the same architecture, hyperparameters, and 10 different random seeds. We are using a batch size of 128, learning rate of 5e-05, and 200 epochs as our default parameter settings. However, in the case that we encounter unexpectedly lower performances on certain tasks, we are planning to further tune our hyperparameters to improve their performance. For random seeds, we are experimenting with 10 different seeds (from 0 to 9) for each task, but will expand to more as time allows.

So far, we’ve tested the T5-small Not-Pretrained model on the SCAN datasets and experimented with various hyperparameters to check if scores are reasonable by referenced documents (Aribandi et al., 2021; Furrer et al., 2020) that are conducting

SCAN (Exact Match)			
seed	simple	length	add-turn-left
0	0.9976	0.1306	0.9528
1	0.9995	0.1306	0.7724
2	0.9864	0.1306	0.9015
3	0.9852	0.1265	0.9462
4	0.9962	0.1306	0.9710
5	0.9758	0.1298	0.7136
6	0.9684	0.1130	0.9636
7	0.9837	0.0957	0.7425
8	0.9840	0.1306	0.8146
9	0.9823	0.1230	0.9553

Figure 1: Exact Match Scores

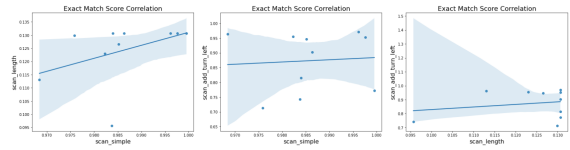


Figure 2: Correlation on SCAN tasks

similar experiments. Tasks we have experimented on the SCAN dataset are “simple”, “length” and “addprim turn left”. After finding the optimal hyperparameters and number of epochs, we have tested out 10 different random seeds. Figure 1 shows our result.

Also, we’ve conducted correlation analysis on three different tasks. Results are all positive but vary from 0.07 to 0.41 depending on the pair of tasks (Figure 2).

Currently, one issue we are facing is in finding the right set of hyperparameters that achieves the expected exact match score. One best hyperparameter setting for one task does not guarantee the best performance of another task. Also, the limitation of gpu allowance is one minor issue that hinders us from expanding our project scope to more datasets and models.

These issues would be solved by starting from hyperparameters that were used from other relevant papers, but needs a sufficient amount of time for us to try numerous settings. Also, we are planning to use the tensorboard to monitor the learning and loss curves to keep track of the performance and select the optimal number of epochs for each model. As we move onto stable architecture for testing, we are planning to submit a batch job to SLURM to allocate the necessary compute resources.

References

- Vamsi Aribandi, Yi Tay, Tal Schuster, Jinfeng Rao, Huaixiu Steven Zheng, Sanket Vaibhav Mehta, Honglei Zhuang, Vinh Q. Tran, Dara Bahri, Jianmo Ni, Jai Gupta, Kai Hui, Sebastian Ruder, and Donald Metzler. 2021. [Ext5: Towards extreme multi-task scaling for transfer learning](#).
- Daniel Furrer, Marc van Zee, Nathan Scales, and Nathanael Schärli. 2020. [Compositional generalization in semantic parsing: Pre-training vs. specialized architectures](#).
- Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. 2019. [Compositionality decomposed: how do neural networks generalise?](#)
- Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. 2019. [Measuring compositional generalization: A comprehensive method on realistic data](#).
- Najoung Kim and Tal Linzen. 2020. [Cogs: A compositional generalization challenge based on semantic interpretation](#).
- Brenden Lake and Marco Baroni. 2018. [Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks](#). In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2873–2882. PMLR.
- Yuanpeng Li, Liang Zhao, Jianyu Wang, and Joel Hestness. 2019. [Compositional generalization for primitive substitutions](#).
- Linqing Liu, Patrick Lewis, Sebastian Riedel, and Pontus Stenetorp. 2021. [Challenges in generalization in open domain question answering](#).
- Pia Weißenhorn, Yuekun Yao, Lucia Donatelli, and Alexander Koller. 2022. [Compositional generalization requires compositional parsers](#).