

BERTs of a feather do not generalize together: Large variability in generalization across models with similar test set performance

R. Thomas McCoy, Junghyun Min, & Tal Linzen

Department of Cognitive Science, Johns Hopkins University

tom.mccoy@jhu.edu, jmin10@jhu.edu, tal.linzen@jhu.edu

Abstract

If the same neural architecture is trained multiple times on the same dataset, will it make similar linguistic generalizations across runs? To study this question, we fine-tuned 100 instances of BERT on the Multi-genre Natural Language Inference (MNLI) dataset and evaluated them on the HANS dataset, which measures syntactic generalization in natural language inference. On the MNLI development set, the behavior of all instances was remarkably consistent, with accuracy ranging between 83.6% and 84.8%. In stark contrast, the same models varied widely in their generalization performance. For example, on the simple case of subject-object swap (e.g., knowing that *the doctor visited the lawyer* does not entail *the lawyer visited the doctor*), accuracy ranged from 0.00% to 66.2%. Such variation likely arises from the presence of many local minima that are equally attractive to a low-bias learner such as a neural network; decreasing the variability may therefore require models with stronger inductive biases.

1 Introduction

Generalization is a crucial component of learning a language. No training set can contain all possible sentences, so learners must be able to generalize to sentences that they have never encountered before. We differentiate two types of generalization:

1. **In-distribution generalization:** Generalization to examples which are novel but which are drawn from the same distribution as the training set.
2. **Out-of-distribution generalization:** Generalization to examples drawn from a different distribution than the training set.

The standard evaluation procedure in natural language processing uses test sets that were generated in the same way as the training set, therefore

testing only in-distribution generalization. Current neural architectures perform very well at this type of generalization. For example, on the natural language understanding tasks in the GLUE benchmark (Wang et al., 2019), several Transformer-based models (Liu et al., 2019b,a; Yang et al., 2019) have surpassed the human baselines established by Nangia and Bowman (2019).

An alternative evaluation approach is to test models on targeted datasets designed to illuminate how the models handle some particular linguistic phenomenon. In this line of investigation, which tests out-of-distribution generalization, the results are more mixed. Some works have found successful handling of subtle linguistic phenomena such as subject-verb agreement (Gulordava et al., 2018) and filler-gap dependencies (Wilcox et al., 2018). Other works, however, have illuminated surprising failures even on seemingly simple types of examples (McCoy et al., 2019). Such results make it clear that there is still much room for improvement in how neural models perform on syntactic structures that are rare in training corpora.

In this work, we investigate whether the linguistic generalization behavior of a given neural architecture is consistent across multiple instances of that architecture. This question is important because, in order to tell which types of architectures generalize best, we need to know whether successes and failures of generalization should be attributed to aspects of the architecture or to random luck in the choice of the model’s initial weights.

We investigate this question using the task of natural language inference (NLI) by fine-tuning 100 instances of BERT (Devlin et al., 2019) on the MNLI dataset (Williams et al., 2018). These 100 instances differed only in (i) the initial weights of the classifier trained on top of BERT and (ii) the order in which examples were sampled during training. All other aspects of training, including

the initial weights of BERT, were held constant. We evaluated these 100 instances on both the in-distribution MNLI development set and also on the out-of-distribution HANS evaluation set (McCoy et al., 2019), which was designed to evaluate syntactic generalization in NLI models.

We find that these 100 instances were remarkably consistent on in-distribution generalization, with all accuracies on the MNLI development set falling in the range 83.6% to 84.8%. In contrast, these 100 instances varied dramatically in their out-of-distribution generalization performance; for example, on one of the thirty categories of examples in the HANS dataset, accuracy ranged from 7% to 77%. These results show that, when assessing the linguistic generalization of neural models, it is important to consider multiple training runs of each architecture, since models can differ vastly in how they perform on examples drawn from a different distribution than the training set, even when they perform similarly on an in-distribution test set.

2 Background

2.1 In-distribution generalization

Several past works have noted that the same architecture can perform very differently across random restarts (Reimers and Gurevych, 2017, 2018). Most relevantly for our work, the original BERT paper (Devlin et al., 2019) noted that fine-tuning of BERT is unstable for some datasets, such that some random restarts achieve state-of-the-art results while others perform poorly; this instability has further been noted in Phang et al. (2018). However, this observation only applies to in-distribution generalization, while we focus on out-of-distribution generalization.

2.2 Out-of-distribution generalization

There have also been several past works that included a focus on how different runs of the same architecture can have different out-of-distribution generalization. Weber et al. (2018) trained 50 instances of a sequence-to-sequence model on a simple symbol replacement task and found that all 50 instances achieved over 99% accuracy on the in-distribution test set but had highly variable performance on out-of-distribution generalization sets; for example, in the most variable case, accuracy ranged from close to 0% to over 90%, with a standard deviation of 20.6%. Similarly, McCoy et al.

(2018) trained 100 instances for each of 6 types of sequence-to-sequence recurrent neural networks, using a synthetic training set that was ambiguous between two generalizations. Some of these models were highly consistent across instances in terms of which generalization they learned, but others varied considerably, with some instances of a given architecture strongly preferring one of the two generalizations while other instances strongly preferred the other generalization. Finally, Liška et al. (2018) trained 5000 instances of recurrent neural networks on the synthetic lookup tables task and found that most failed on compositional generalization but that a small number did display strong compositional generalization.

All of these past works that studied variation in out-of-distribution generalization used simple, synthetic tasks with training sets carefully designed to have certain types of examples withheld. Our work extends this area of inquiry to models trained on natural language, which is noisier and does not have such explicit constraints, to see if models are still as variable even in more practical cases where the training set is not adversarially designed to be ambiguous.

2.3 Linguistic analysis of BERT

Many recent papers have sought a deeper understanding of BERT, whether to assess its encoding of sentence structure (Lin et al., 2019; Hewitt and Manning, 2019; Chrupała and Alishahi, 2019; Jawahar et al., 2019; Tenney et al., 2019b); its representational structure more generally (Abnar et al., 2019); its handling of specific linguistic phenomena such as subject-verb agreement (Goldberg, 2019), negative polarity items (Warstadt et al., 2019), function words (Kim et al., 2019), or a variety of psycholinguistic phenomena (Ettinger, 2019); or its internal workings (Coenen et al., 2019; Tenney et al., 2019a; Clark et al., 2019). The novel contribution of this work is the focus on variability across a large number of fine-tuning runs; previous works have generally used models without fine-tuning or have used only a small number of fine-tuning runs (usually only one fine-tuning run, or at most five fine-tuning runs).

3 Method

3.1 Task and datasets

We used the task of natural language inference (NLI; also known as Recognizing Textual Entail-

Heuristic	Definition	Example
Lexical overlap	Assume that a premise entails all hypotheses constructed from words in the premise	The doctor was paid by the actor. $\xrightarrow{\text{WRONG}}$ The doctor paid the actor.
Subsequence	Assume that a premise entails all of its contiguous subsequences.	The doctor near the actor danced. $\xrightarrow{\text{WRONG}}$ The actor danced.
Constituent	Assume that a premise entails all complete subtrees in its parse tree.	If the artist slept, the actor ran. $\xrightarrow{\text{WRONG}}$ The artist slept.

Figure 1: The heuristics targeted by the HANS dataset, along with examples of incorrect entailment predictions that these heuristics would lead to. (Figure from McCoy et al. (2019).)

ment), which involves giving a model two sentences, called the *premise* and the *hypothesis*. The model must then output a label of *entailment* if the premise entails (i.e., implies the truth of) the hypothesis, *contradiction* if the premise contradicts the hypothesis, or *neutral* otherwise. For training, we used the training set of the MNLI dataset (Williams et al., 2018), examples from which are given below:

- (1)
 - a. **Premise:** One of our number will carry out your instructions minutely.
 - b. **Hypothesis:** A member of my team will execute your orders with immense precision.
 - c. **Label:** entailment
- (2)
 - a. **Premise:** but that takes too much planning
 - b. **Hypothesis:** It doesn't take much planning.
 - c. **Label:** contradiction
- (3)
 - a. **Premise:** He turned and smiled at Vrenna.
 - b. **Hypothesis:** He smiled at Vrenna who was walking slowly behind him with her mother.
 - c. **Label:** neutral

To test models' in-distribution generalization, we evaluated their performance on the MNLI matched development set, which was generated in the same way as the MNLI training set.¹ To test out-of-distribution generalization, we used the HANS dataset (McCoy et al., 2019), which con-

¹We used the development set rather than the test set because the test set labels are not available to the public. This development set was not used in any way during training, making it effectively a test set.

tains NLI examples designed to require understanding of syntactic structure. More specifically, the HANS dataset targets three structural heuristics that models trained on MNLI are likely to learn. These heuristics are defined in Figure 1, along with examples of cases where the heuristics make incorrect predictions.

To assess whether a model has learned these heuristics, the HANS dataset contains examples where each heuristic makes the right predictions (i.e., where the correct label is *entailment*) and examples where each heuristic makes the wrong predictions (i.e., where the correct label is *non-entailment*). A model that has adopted one of the three heuristics will output *entailment* for all examples targeting that heuristic, even the examples for which the correct answer is *non-entailment*.

3.2 Models and training

All of our models consisted of BERT with a linear classifier on top of it outputting labels of *entailment*, *contradiction*, or *neutral*. We fine-tuned this model on MNLI using the MNLI fine-tuning code available on the BERT GitHub repository.² We initialized the BERT component of the model with the pre-trained `bert-base-uncased` weights from the BERT GitHub repository; these weights were obtained by training BERT on a large quantity of natural text. All of our instances of fine-tuning used these same initial BERT weights, but they all used different, random initial weights for the classifier. The fine-tuning process then modified the weights of both the BERT component and the classifier. We ran 100 instances of this fine-tuning. Following Devlin et al. (2019), we varied only two things across fine-tuning runs: (i)

²<https://github.com/google-research/bert>

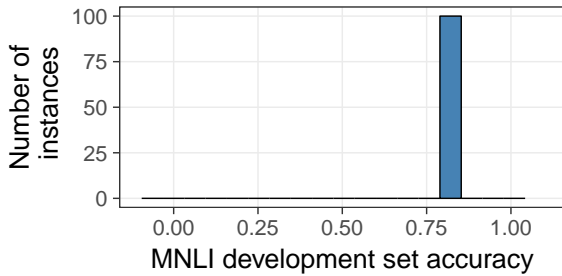


Figure 2: In-distribution generalization: Performance on the MNLI development set

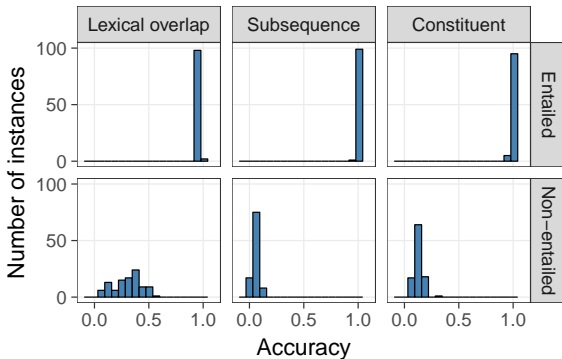


Figure 3: Out-of-distribution generalization: Performance on the HANS evaluation set, broken down into six categories of examples based on which syntactic heuristic each example targets and whether the correct label is *entailment* or *non-entailment*. The non-entailed lexical overlap cases (lower left plot) display a large degree of variability across instances.

the random initial weights of the classifier and (ii) the order in which examples were sampled during training. All other aspects, including the initial pre-trained weights of the BERT component, were held constant across all 100 runs.

4 Results

4.1 In-distribution generalization

The 100 instances performed remarkably consistently on in-distribution generalization, with all models achieving an accuracy between 83.6% and 84.8% on the MNLI development set (Figure 2).

4.2 Out-of-distribution generalization

On the HANS evaluation set, performance was much more variable. This evaluation set consists of 6 main categories of examples, each of which can be further divided into 5 subcategories. Performance was reasonably consistent on five of the six categories, but on the sixth category—the non-entailed lexical overlap category—performance

varied dramatically, ranging from 6% accuracy to 54% accuracy (Figure 3). Given that this is the most variable case, we focus on it for the rest of the analysis of the HANS results.

The non-entailed lexical overlap category encompasses examples for which the correct label is non-entailment (i.e., either *contradiction* or *neutral*) and for which all the words in the hypothesis also appear in the premise but not as a contiguous subsequence of the premise. This category can be broken down into five subcategories; examples for these subcategories, along with the distribution of the instances’ performance on each subcategory, are in Figure 4. Chance performance on HANS is 50%; on all of the subcategories except for passive sentences, accuracies range from far below chance to modestly above chance. At least some of these subcategories seem intuitively simple, yet models still vary considerably on them; for example, the subject-object swap examples could be handled with only rudimentary knowledge of syntax (in particular, knowledge of the distinction between subjects and objects), yet models still range in accuracy on this subcategory from 0% to 66%, indicating that, although these models performed very consistently on the in-distribution test set, they have nevertheless learned strikingly variable representations of syntax.

5 Discussion

We have found that models that vary only in their initial weights and the order of training examples can vary substantially in their out-of-distribution linguistic generalization. We found this variation even with the vast majority of initial weights held constant (i.e., all the weights in the BERT component of the model), suggesting that models might display an even greater degree of variability if the pre-training step used to initialize the weights of the BERT component were also redone across instances. These results underscore the importance of evaluating models on multiple random restarts, as conclusions drawn from a single instance of a model might not hold across instances. Further, these results also highlight the importance of evaluating models on out-of-distribution generalization; given that all 100 of our instances displayed similar in-distribution generalization, only their performance on out-of-distribution generalization actually illuminates the substantial differences in what these models have learned.

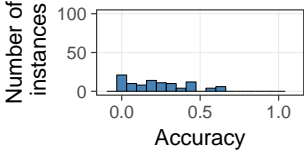
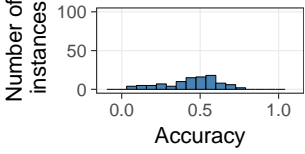
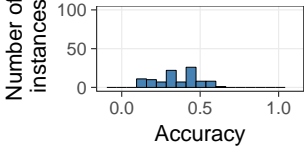
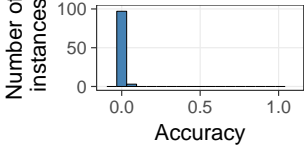
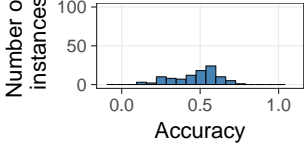
Category	Example	Accuracy distribution
Subject-object swap	The doctor visited the lawyer. → The lawyer visited the doctor.	
Preposition	The judges by the manager saw the artists. → The artists saw the manager.	
Relative clause	The actors advised the author who the tourists saw. → The author advised the tourists.	
Passive	The senators were recommended by the managers. → The senators recommended the managers.	
Conjunction	The doctors advised the presidents and the tourists. → The presidents advised the tourists.	

Figure 4: Results on the various subcategories within the non-entailed lexical overlap examples of the HANS dataset. We do not include the other 25 subcategories of the HANS dataset in this figure as there was little variability across instances for those subcategories.

The high degree of variability shown by these models likely reflects the presence of many local minima in the loss surface, all of which are equally attractive to our models, making the choice of the particular minimum that the model settles on essentially arbitrary and easily affected by random variations in initial weights and the order of training examples. To reduce this variability, therefore, it will likely be necessary to use models with stronger inductive biases that can help distinguish between these many local minima. In stark contrast to the models we have looked at—which generalized in highly variable ways despite being trained on the same set of examples—humans tend to converge to remarkably similar linguistic generalizations despite major differences in the linguistic input that they encounter as children (Chomsky, 1965, 1980). This fact suggests that humans have stronger inductive biases than these models, leading to more robustly similar generalization in humans than we have observed with our re-runs of BERT. This suggests that reducing the general-

ization variability of NLP models is desirable as a step toward bringing models closer to human performance in one of the major areas where they still dramatically lag behind humans, namely in out-of-distribution generalization.

Acknowledgments

We are grateful to Emily Pitler, Dipanjan Das, and the members of the Johns Hopkins Computation and Psycholinguistics lab group for helpful comments. Any errors are our own.

This project is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. 1746891 and by a gift from Google, and it was conducted using computational resources from the Maryland Advanced Research Computing Center (MARCC). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation, Google, or MARCC.

References

- Samira Abnar, Lisa Beinborn, Rochelle Choenni, and Willem Zuidema. 2019. [Blackbox meets blackbox: Representational similarity & stability analysis of neural language models and brains](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 191–203, Florence, Italy. Association for Computational Linguistics.
- Noam Chomsky. 1965. *Aspects of the Theory of Syntax*. MIT Press, Cambridge, MA.
- Noam Chomsky. 1980. Rules and representations. *Behavioral and Brain Sciences*, 3(1):1–15.
- Grzegorz Chrupała and Afra Alishahi. 2019. [Correlating neural and symbolic representations of language](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2952–2962, Florence, Italy. Association for Computational Linguistics.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? An analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. 2019. [Visualizing and measuring the geometry of BERT](#). *arXiv preprint arXiv:1906.02715*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Allyson Ettinger. 2019. [What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models](#). *arXiv preprint arXiv:1907.13528*.
- Yoav Goldberg. 2019. [Assessing BERT’s syntactic abilities](#). *arXiv preprint arXiv:1901.05287*.
- Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. [Colorless green recurrent networks dream hierarchically](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1195–1205, New Orleans, Louisiana. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Najoung Kim, Roma Patel, Adam Poliak, Patrick Xia, Alex Wang, Tom McCoy, Ian Tenney, Alexis Ross, Tal Linzen, Benjamin Van Durme, Samuel R. Bowman, and Ellie Pavlick. 2019. [Probing what different NLP tasks teach machines about function word comprehension](#). In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 235–249, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. [Open sesame: Getting inside BERT’s linguistic knowledge](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 241–253, Florence, Italy. Association for Computational Linguistics.
- Adam Liška, Germán Kruszewski, and Marco Baroni. 2018. [Memorize or generalize? searching for a compositional rnn in a haystack](#). In *Proceedings of the 2018 workshop on Architectures and Evaluation for Generality, Autonomy, and Progress in AI (AEGAP)*.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019a. [Multi-task deep neural networks for natural language understanding](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4487–4496, Florence, Italy. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. [RoBERTa: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- R. Thomas McCoy, Robert Frank, and Tal Linzen. 2018. [Revisiting the poverty of the stimulus: Hierarchical generalization without a hierarchical bias in recurrent neural networks](#). In *Proceedings of the 40th Annual Conference of the Cognitive Science Society*, pages 2093–2098, Madison, WI.
- R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the wrong reasons: Diagnosing](#)

- syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.
- Nikita Nangia and Samuel R. Bowman. 2019. **Human vs. muppet: A conservative estimate of human performance on the GLUE benchmark**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4566–4575, Florence, Italy. Association for Computational Linguistics.
- Jason Phang, Thibault Févry, and Samuel R Bowman. 2018. **Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks**. *arXiv preprint arXiv:1811.01088*.
- Nils Reimers and Iryna Gurevych. 2017. **Reporting score distributions makes a difference: Performance study of LSTM-networks for sequence tagging**. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 338–348, Copenhagen, Denmark. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2018. **Why comparing single performance scores does not allow to draw conclusions about machine learning approaches**. *arXiv preprint arXiv:1803.09578*.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019a. **BERT rediscovers the classical NLP pipeline**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019b. **What do you learn from context? probing for sentence structure in contextualized word representations**. In *International Conference on Learning Representations*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. **GLUE: A multi-task benchmark and analysis platform for natural language understanding**. In *International Conference on Learning Representations*.
- Alex Warstadt, Yu Cao, Ioana Grosu, Wei Peng, Hagen Blix, Yining Nie, Anna Alsop, Shikha Bordia, Haokun Liu, Alicia Parrish, Sheng-Fu Wang, Jason Phang, Anhad Mohananey, Phu Mon Htut, Paloma Jeretic, and Samuel R. Bowman. 2019. **Investigating BERT’s knowledge of language: Five analysis methods with NPIs**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2870–2880, Hong Kong, China. Association for Computational Linguistics.
- Noah Weber, Leena Shekhar, and Niranjan Balasubramanian. 2018. **The fine line between linguistic generalization and failure in Seq2Seq-attention models**. In *Proceedings of the Workshop on Generalization in the Age of Deep Learning*, pages 24–27, New Orleans, Louisiana. Association for Computational Linguistics.
- Ethan Wilcox, Roger Levy, Takashi Morita, and Richard Futrell. 2018. **What do RNN language models learn about filler-gap dependencies?** In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 211–221, Brussels, Belgium. Association for Computational Linguistics.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. **A broad-coverage challenge corpus for sentence understanding through inference**. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. 2019. **XLNet: Generalized autoregressive pretraining for language understanding**. *arXiv preprint arXiv:1906.08237*.