

Mediators in Determining what Processing BERT Performs First

Aviv Slobodkin

Leshem Choshen

Omri Abend

School of Computer Science and Engineering

The Hebrew University of Jerusalem

{aviv.slobodkin, leshem.choshen, omri.abend}@mail.huji.ac.il

Abstract

Probing neural models for the ability to perform downstream tasks using their activation patterns is often used to localize what parts of the network specialize in performing what tasks. However, little work addressed potential mediating factors in such comparisons. As a test-case mediating factor, we consider the prediction’s *context length*, namely the length of the span whose processing is minimally required to perform the prediction. We show that not controlling for context length may lead to contradictory conclusions as to the localization patterns of the network, depending on the distribution of the probing dataset. Indeed, when probing BERT with seven tasks, we find that it is possible to get 196 different rankings between them when manipulating the distribution of context lengths in the probing dataset. We conclude by presenting best practices for conducting such comparisons in the future.¹

1 Introduction

The strong performance of end-to-end models and the difficulty in understanding their inner workings has led to extensive research aimed at interpreting their behavior (Li et al., 2016; Yosinski et al., 2015; Karpathy et al., 2015). This notion has led researchers to investigate the behavioral traits of networks in general (Li et al., 2015; Hachon et al., 2020) and representative architectures in particular (Schlichtkrull et al., 2020). Within NLP, Transformer-based pretrained embeddings are the basis for many tasks, which underscores the importance in interpreting their behavior (Belinkov et al., 2020), and especially the behavior of BERT (Devlin et al., 2019; Rogers et al., 2020), perhaps the most widely used of Transformer-based models.

In this work, we analyze the common approach of *probing* (§2), used to localize where “knowledge”

of particular tasks is encoded; localization is often carried out in terms of the layers most responsible for the task at hand (c.f. Tenney et al., 2019b). Various works (Tenney et al., 2019a; Peters et al., 2018; Blevins et al., 2018) showed that some tasks are processed in lower levels than others.

We examine the extent to which potential mediating factors may account for observed trends and show that varying some mediating factors (see §2) may diminish, or even reverse, the conclusions made by Tenney et al. (T19; 2019a). Specifically, despite reaffirming T19’s experimental findings, we contest T19’s interpretation of the results, namely that the processing carried out by BERT parallels the classical NLP pipeline. Indeed, T19 concludes that lexical tasks (POS tagging) are performed by the lower layers, followed by syntactic tasks, whereas more semantic tasks are performed later on. This analysis rests on the assumption that the nature of the task (lexical, syntactic, or semantic) is the driving force that determines what layer performs what analysis. We show that other factors should be weighed in as well. Specifically, we show that manipulating the distribution of examples in the probing dataset can lead to a variety of different conclusions as to what tasks are performed first.

We argue that potential mediators must be considered when comparing tasks, and focus on one such mediator – the *context length*, which we define as the number of tokens whose processing is minimally required to perform the prediction. We operationalize this notion by defining it as the maximal distance between any two tokens for which a label is predicted. This amounts to the span length in tasks that involve a single span (e.g., NER), and to the dependency length in tasks that address the relation between two spans. See §2. Our motivation for considering context length as a mediator is grounded in previous work that presented the difficulty posed by long-distance dependencies in various NLP tasks (Xu et al., 2009; Sennrich, 2017),

¹The code is available at <https://github.com/lovodkin93/BERT-context-distance>.

and particularly in previous work that indicated the Transformers’ difficulty to generalize across different dependency lengths (Choshen and Abend, 2019).

We show that in some of the cases where one task seems to be better predicted by a higher layer than another task, controlling for context length may reverse that order. Indeed we show that 196 different rankings between the seven tasks explored in T19 may be obtained with a suitable distribution over the probing datasets, namely 196 different ways to rank the tasks according to their expected layer. Moreover, our results show that when context length is not taken into account, one task (e.g., dependency parsing) may seem to be processed at a higher layer than another (e.g., NER), when its expected layer (see §2) is, in fact, lower for all ranges of context lengths (§3.1.1).

2 Background

We begin by laying out the terminology and methodology we will use in the paper.

Edge Probing. Edge probing is the method of training a classifier for a given task on different parts of the network (without fine-tuning). Success in classification is interpreted as evidence that the required features for classification are somehow encoded in the examined part and are sufficiently easy to extract. In our experiments, we follow T19 and probe BERT with Named Entity Recognition (NER), a constituent-based task (classifying Non-terminals - Non-term.), Semantic Role Labeling (SRL), Co-reference (Co-ref.), Semantic Proto-Roles (SPR; Reisinger et al., 2015), Relation Classification (RC) and the Stanford Dependency Parsing (Dep.; de Marneffe et al., 2006).

Causal considerations in interpreting probing results were also emphasized by several recent works (e.g., Kaushik et al., 2020; Vig et al., 2020; Elazar et al., 2021).

Localization by Expected Layer. The expected layer metric (which we will henceforth refer to it as \mathbb{E}_{layer}) of T19 assesses which layer in BERT is most needed for prediction: a probing classifier $P^{(l)}$ is trained on the lowest l layers. Then, a differential score $\Delta^{(l)}$ is computed, which indicates the performance gain when taking into account one additional layer:

$$\Delta^{(l)} = \text{Score}(P^{(l)}) - \text{Score}(P^{(l-1)}) \quad (1)$$

Once all the $\{\Delta^{(l)}\}_{l=1}^{12}$ are computed, we may compute \mathbb{E}_{layer} :

$$\mathbb{E}_{layer}[l] = \frac{\sum_{l=1}^{12} l \cdot \Delta^{(l)}}{\sum_{l=1}^{12} \Delta^{(l)}} \quad (2)$$

Therefore, unlike standard edge probing, which is performed on each layer individually, computing \mathbb{E}_{layer} takes into account all layers up to a given l .

Mediation Analysis. Each of the explored tasks classifies one or two input sub-spans. In both cases, we define the context length to be the distance between the earliest and latest span index. Namely, for tasks with two spans (e.g., SPR), $span_1=[i_1, j_1]$ and $span_2=[i_2, j_2]$, where $span_1$ appears before $span_2$, the context length is $j_2 - i_1$, whereas for tasks with just one span (e.g., NER), $span_1=[i_1, j_1]$, it is $j_1 - i_1$.

In order to examine the effect of context length on \mathbb{E}_{layer} , we model it as a mediating factor, namely as an intermediate variable that (partly) explains the relationship between two other variables (in this work, a task and its \mathbb{E}_{layer}). See Figure 1.

We bin each task’s test set into non-overlapping bins, according to their context length ranges. We use the notation ‘i-j’ to denote the bin of context lengths in the range [i,j]. For example, the second bin would be ‘3-5’, denoting context lengths 3, 4, and 5. In addition, given a specific task, two possible approaches exist to examine the mediation effect of context length on the task’s \mathbb{E}_{layer} . The first one bins all the task’s data into sub-sets, in advance. Then, this approach fine-tunes over each subset separately. Alternatively, the second approach fine-tunes over the whole dataset, binning only during the test phase. We follow the latter approach, as it is more computationally efficient.

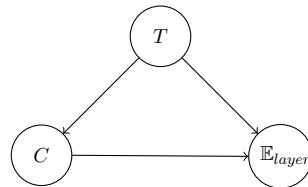


Figure 1: The relationship we stipulate between the task, the context length, and \mathbb{E}_{layer} . We use two random variables: T is the task, which can be any of the seven tasks we observe and C is the context length.

Interestingly, in §3.1.1, we encounter a special edge case, where the aggregated average (i.e., \mathbb{E}_{layer}) of one task is higher than another, whereas

in each sub-set (by a given context length) it is lower. This may occur when the weight of the sub-sets differs between the two aggregations.

3 Experiments

We hypothesize that the context length is a mediating factor in the \mathbb{E}_{layer} of a task. In order to test this hypothesis, we run the following experiments, aiming at isolating the context length.

We use the SPR1 dataset (Reisinger et al., 2015) to probe SPR, the English Web Treebank for the Dep. task (Silveira et al., 2014), the SemEval 2010 Task 8 for the RC task (Hendrickx et al., 2009), and the OntoNotes 5.0 dataset (Weischedel et al., 2013) for the other tasks. Configurations follow the defaults in the Jiant toolkit implementation (Wang et al., 2019). In addition, we work with the BERT-base model.

3.1 The Effect on \mathbb{E}_{layer}

First, we wish to confirm that context length indeed affects \mathbb{E}_{layer} and that the task is not a sole contributor to this. Given a task and a threshold thr , we compile a dataset for the task containing the subset of examples with context lengths shorter than thr , and use it to compute \mathbb{E}_{layer} . We do it for all tasks and for every integer threshold between 0 and a maximal threshold, which is selected separately for each task to ensure that at least 2000 instances remain in the last bin.

We find that context length plays an important role in the difference between the expected layers (Figure 2). Most notably, the Co-ref., SRL, Dep., and RC tasks’ \mathbb{E}_{layer} increases when increasing the threshold.

Next, we divide the data into smaller bins of non-overlapping context length ranges, in order to control for the influence of the context lengths on the expected layers of the tasks. We compute \mathbb{E}_{layer} for sub-sets of similar lengths. In choosing the size of each such range, we try to balance between informativeness (narrower ranges) and reliability (having enough examples in each range, so as to reduce noise). We find that the narrowest range width that retains at least 1% of the examples in each bin is 3. We thus divide the dataset for each task into context length ranges of width 3, until the maximal threshold is reached. Higher context lengths are lumped into an additional bin.

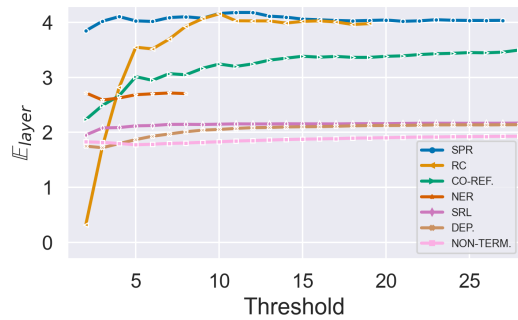


Figure 2: \mathbb{E}_{layer} as a function of a threshold on the context length. For each such threshold thr (x-axis), \mathbb{E}_{layer} (y-axis) is computed based only on the examples with context length no longer than thr .

3.1.1 Manipulating the Context Length Distribution: An Extreme Case.

We begin by examining two specific tasks: Dep. and NER, and their \mathbb{E}_{layer} for each context length’s range. We then consider, for simplicity, a case where all the context lengths of Dep. are of length 9+, while those of NER are in the range of 3-5 (Figure 3). We see that when controlling for context length, Dep. is computed in a lower layer than NER, regardless of the range. However, depending on the distribution of context lengths in the probing dataset, the outcome may be completely different, with Dep. being processed in higher layers (for a similar example of a different task-pair, see §A.1).

These results indicate that the results of T19 do not necessarily indicate that BERT is performing a pipeline of computations (as is commonly asserted, see e.g., T19 and Blevins et al. (2018)), and that mediating factors need to be taken into account when interpreting \mathbb{E}_{layer} .

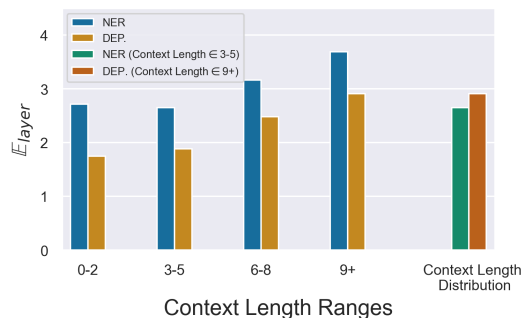


Figure 3: \mathbb{E}_{layer} of NER and Dep. for different context length ranges (4 left blue and yellow pairs), and their \mathbb{E}_{layer} when all instances of NER are of context length $l \in [3, 5]$ and all those of Dep. are of context length $l \geq 9$ (rightmost green and red pair). While for every context length range, NER’s \mathbb{E}_{layer} is bigger than that of Dep., for some context length distribution that order may be reversed.

3.2 Imposing Similar Length Distributions

In the previous section, we observed that one task can be both higher and lower than another. That depends on the distribution of context lengths in the probing dataset. We next ask whether such a "paradox" arises in experiments when imposing the same context length distributions on the two tasks.

Following Pearl (2001), we employ mediation analysis and specifically concentrate on the Natural Direct Effect (NDE), which is the difference between two of the observed dependent variables (in our case \mathbb{E}_{layer}), when fixing the mediator. In our case, the NDE is the difference between the \mathbb{E}_{layer} of two tasks, while forcing the same context length distribution on both. For convenience, we force the distribution of one of the examined tasks (for more details, see §A.2), but any distribution is applicable. In general, the equation for computing the NDE of tasks t_1 and t_2 , with the context length distribution of t_1 imposed on both, is:

$$NDE_{t_1 \rightarrow t_2} = \sum_c [E_{\Delta}[l|C=c, T=t_2] - E_{\Delta}[l|C=c, T=t_1]] \cdot P(C=c|T=t_1) \quad (3)$$

where T is a random variable of the tasks, and C is a random variable of the context length.

We apply NDE twice for every pair of tasks (once for each task’s context length distribution). We then compare the results to the difference between the tasks’ expected layers where each task keeps its original context length distribution (unmediated). Results (Figure 4) show that the difference could be more than 50 times larger (change of 1.24 in absolute value) or decrease by 86% (0.73 in absolute value). In some cases the order of the two tasks is reversed, namely, the task that is lower with one distribution becomes higher with another. This shows that even among our examined set of seven tasks, the effect of potential mediators cannot be ignored. For more results, see §A.3.

3.2.1 Controlling for Context Length

After observing that the distribution of context length in the probing dataset may affect the relative order of the expected layers, we propose a more detailed and accurate method to compare the expected layers, which does not rely on a specific length distribution. We do so by plotting the *controlled effect*, namely \mathbb{E}_{layer} for each range separately.

Our results (Figure 5) allow computing the range of possible expected layers for a task, that may result from taking any context length distribution

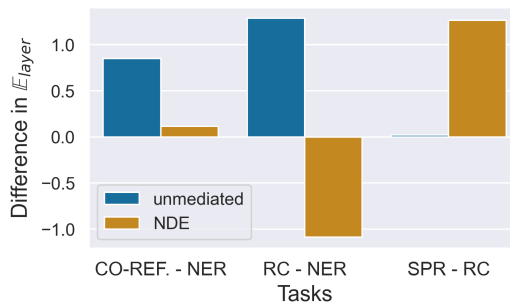


Figure 4: Difference between unmediated \mathbb{E}_{layer} and NDE for NER and Co-ref. (left); NER and RC (middle); and SPR and RC (right). The employed context length distributions (as part of the NDE calculations) are of Co-ref., NER and SPR, respectively.

(Figure 6). The figure shows the wide range of possible relative behaviors of \mathbb{E}_{layer} for task-pairs: from notable to negligible difference in expected layers (e.g., SRL and Co-ref.), to pairs whose ordering of expected layers may be reversed (i.e., overlapping ranges, such as with SPR and RC). In fact, by taking into account every possible combination of context length distribution for each of the tasks, we get as many as 196 possible rankings of the seven tasks according to their \mathbb{E}_{layer} . One such possible order is, for example, Non-term. < Dep. < SRL < RC < NER < Co-ref. < SPR. We elaborate on this in §A.4.

To recap, we find that the difference in \mathbb{E}_{layer} between some tasks may considerably change and their order may reverse, depending on the context length. This finding lends further support to our claim that mediators should be taken into account.

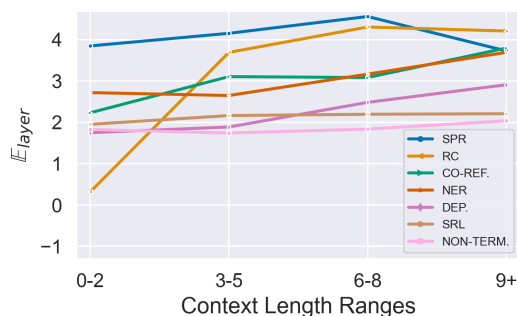


Figure 5: Expected layers of all seven tasks as a function of context length range.

4 Conclusion

We showed that when performing edge probing to identify what layers are responsible for addressing what tasks, it is imperative to take into account potential mediators, as they may be responsible

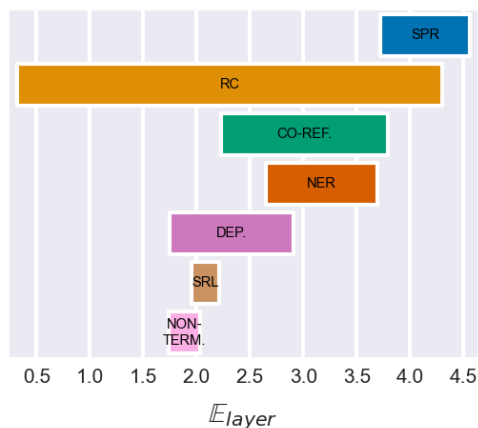


Figure 6: The range of possible expected layers when varying context length, for each of the seven tasks.

for much of the observed effect. Specifically, we showed that context length has a significant impact on a task’s \mathbb{E}_{layer} . Our analysis shows the wide range of relative orderings of the expected layers for different tasks when assuming different context length distributions; from extreme edge cases, like the one we observed in §3.1.1, to more common, but potentially misleading ones, where the difference between expected layers may dramatically increase or decrease depending on the context length distribution. Most importantly, it shows that by manipulating the context length distribution, we may get a wide range of outcomes.

Our work suggests that mediating factors should be taken into account when basing analysis on the \mathbb{E}_{layer} . On a broader note, alternative hypotheses should be considered, before limiting oneself to a single interpretation.

Future work will consider the effect of other mediating factors. The two methods we used, NDE and controlled effect, can be used to examine the impact of other mediating factors and should be adopted as part of the field’s basic analysis toolkit (cf. Feder et al., 2020; Vig et al., 2020). NDE should be used when several effects are examined simultaneously, as it facilitates the assessment of their effect on the tasks’ complexity. It is also advisable to use NDE when a more practical examination is required, i.e., when distributions of the mediators are given empirically, as it is easier to derive the mediating factors’ impact using this method. In contrast, the controlled effect method should be used when examining the effects of two variables (e.g., tasks and mediating factors) or when comparing several tasks with one mediating effect.

Acknowledgements

This work was supported by the Israel Science Foundation (grant no. 929/17). We would also like to thank Amir Feder for his very insightful feedback on our paper.

References

- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2020. [On the linguistic representational power of neural machine translation models](#). *Computational Linguistics*, 46(1):1–52.
- Terra Blevins, Omer Levy, and Luke Zettlemoyer. 2018. [Deep RNNs encode soft hierarchical syntax](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 14–19, Melbourne, Australia. Association for Computational Linguistics.
- Leshem Choshen and Omri Abend. 2019. [Automatically extracting challenge sets for non-local phenomena in neural machine translation](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 291–303, Hong Kong, China. Association for Computational Linguistics.
- Marie-Catherine de Marneffe, Bill MacCartney, and Christopher D. Manning. 2006. [Generating typed dependency parses from phrase structure parses](#). In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC’06)*, Genoa, Italy. European Language Resources Association (ELRA).
- 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.
- Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. [Amnesic probing: Behavioral explanation with amnesic counterfactuals](#). *Transactions of the Association for Computational Linguistics*, 9:160–175.
- Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2020. [Causalm: Causal model explanation through counterfactual language models](#).
- Guy Hacohen, Leshem Choshen, and D. Weinshall. 2020. [Let’s agree to agree: Neural networks share classification order on real datasets](#). *International Conference of Machine Learning*.

- Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2009. [SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals](#). In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions (SEW-2009)*, pages 94–99, Boulder, Colorado. Association for Computational Linguistics.
- Andrej Karpathy, Justin Johnson, and Li Fei-Fei. 2015. [Visualizing and understanding recurrent networks](#).
- Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. 2020. [Learning the difference that makes a difference with counterfactually-augmented data](#). In *International Conference on Learning Representations*.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. [Visualizing and understanding neural models in NLP](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 681–691, San Diego, California. Association for Computational Linguistics.
- Y. Li, J. Yosinski, J. Clune, H. Lipson, and J. Hopcroft. 2015. Convergent learning: Do different neural networks learn the same representations? In *FE@NIPS*.
- Judea Pearl. 2001. Direct and indirect effects. In *Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, UAI '01*, page 411–420, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Matthew Peters, Mark Neumann, Luke Zettlemoyer, and Wen-tau Yih. 2018. [Dissecting contextual word embeddings: Architecture and representation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1499–1509, Brussels, Belgium. Association for Computational Linguistics.
- Drew Reisinger, Rachel Rudinger, Francis Ferraro, Craig Harman, Kyle Rawlins, and Benjamin Van Durme. 2015. [Semantic proto-roles](#). *Transactions of the Association for Computational Linguistics*, 3:475–488.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. [A primer in bertology: What we know about how bert works](#).
- M. Schlichtkrull, Nicola De Cao, and Ivan Titov. 2020. Interpreting graph neural networks for nlp with differentiable edge masking.
- Rico Sennrich. 2017. [How grammatical is character-level neural machine translation? assessing MT quality with contrastive translation pairs](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics.
- Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. 2014. [A gold standard dependency corpus for English](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).
- 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*.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. [Causal mediation analysis for interpreting neural nlp: The case of gender bias](#).
- Alex Wang, Ian F. Tenney, Yada Pruksachatkun, Katherin Yu, Jan Hula, Patrick Xia, Raghu Pappagari, Shuning Jin, R. Thomas McCoy, Roma Patel, Yinghui Huang, Jason Phang, Edouard Grave, Najoung Kim, Phu Mon Htut, Thibault F'evry, Berlin Chen, Nikita Nangia, Haokun Liu, , Anhad Mohananey, Shikha Bordia, Ellie Pavlick, and Samuel R. Bowman. 2019. [jiant 1.0: A software toolkit for research on general-purpose text understanding models](#). <http://jiant.info/>.
- Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2013. Ontonotes release 5.0 ldc2013t19. *Linguistic Data Consortium, Philadelphia, PA*, 23.
- Peng Xu, Jaeho Kang, Michael Ringgaard, and Franz Och. 2009. Using a dependency parser to improve smt for subject-object-verb languages. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, NAACL '09*, page 245–253, USA. Association for Computational Linguistics.
- Jason Yosinski, Jeff Clune, Anh Mai Nguyen, Thomas J. Fuchs, and Hod Lipson. 2015. [Understanding neural networks through deep visualization](#). *CoRR*, abs/1506.06579.

A Appendix

A.1 Additional Example of the Extreme Case

We show another example of a task-pair that, under certain distributions of context lengths, exhibits similar behavior to that observed in the edge case described in §3.1.1 (figure 7).

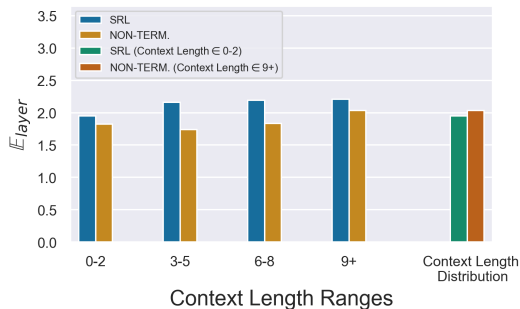


Figure 7: \mathbb{E}_{layer} of SRL and Non-term. for different context length ranges (4 left blue and yellow pairs), and their \mathbb{E}_{layer} when all instances of SRL are of context length $l \in [0, 2]$ and all those of Non-term. are of context length $l \geq 9$ (rightmost green and red pair). While for every context length range, SRL’s \mathbb{E}_{layer} is bigger than that of Non-term., for some context length distribution that order may be reversed.

A.2 Context Length Distribution

A lot of our work deals with possible context length distributions, normalizing distribution, and accounting for the distribution. We provide here the actual distributions which are the underlying property controlling the seen effects. We provide data on the percentage of examples in each context length range for each task (figure 8).

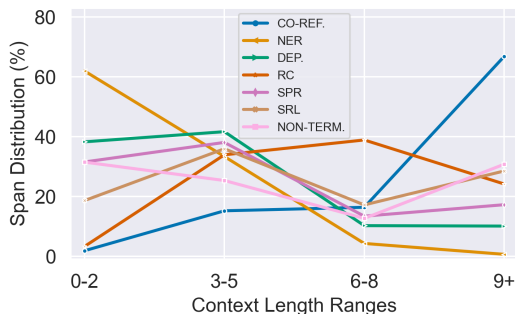


Figure 8: Percentage of examples as a function of context length range, for each of the 7 tasks (see legend).

A.3 NDE vs. Unmediated Difference for All Task-Pairs

For every task-pair, we compare the unmediated \mathbb{E}_{layer} difference with the pair’s NDE. Figure 9 presents this comparison for each task-pair, with

the distribution of one of the pair’s tasks being applied in the NDE calculations, for each task-pair.

A.4 Extreme \mathbb{E}_{layer} Differences

Based on figure 6, we compute the extreme \mathbb{E}_{layer} differences of each task-pair. Namely, for each such pair, we juxtapose the difference between the maximal possible \mathbb{E}_{layer} of the first task and the minimal \mathbb{E}_{layer} of the second one with the opposite case (the difference between the minimal possible \mathbb{E}_{layer} of the first task and the maximal \mathbb{E}_{layer} of the second one). Our results can be seen in figure 10.

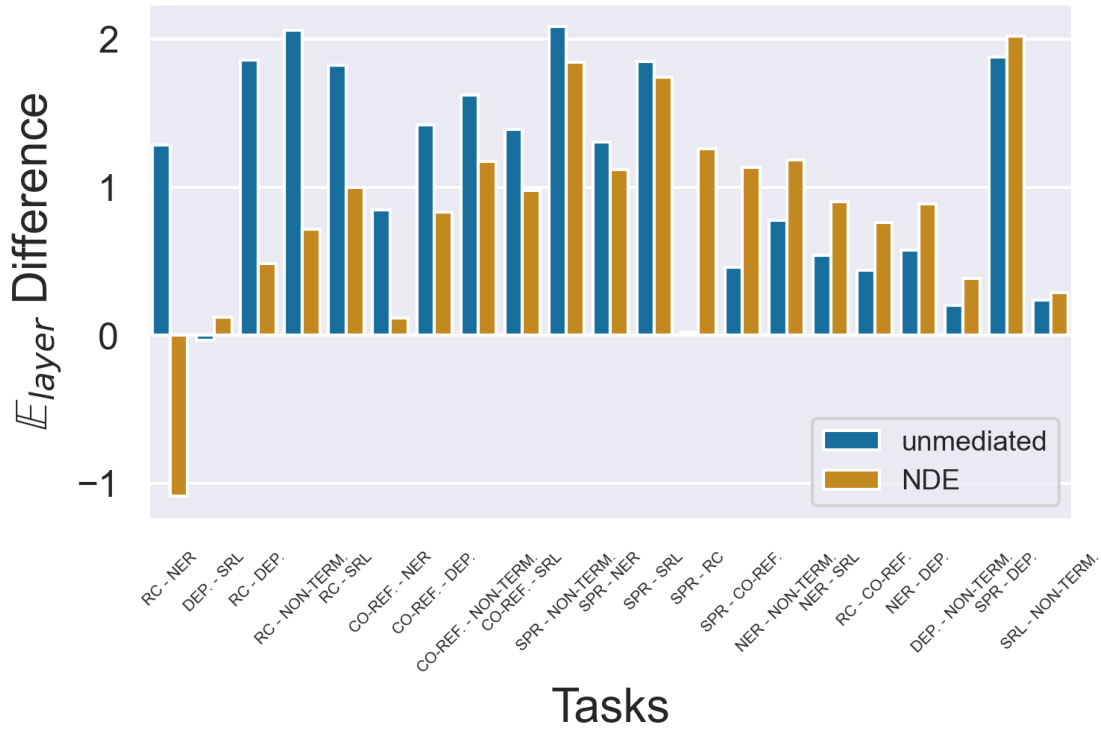


Figure 9: Difference between unmediated \mathbb{E}_{layer} and NDE for every task-pair. The employed context length distributions (as part of the NDE calculations) are, from left to right, of NER, SRL, Dep., Non-term., SRL, Co-ref., Dep., Non-term., SRL, Non-term., SPR, SRL, SPR, SPR, Non-term., SRL, RC, NER, Non-term., Dep. and SRL.

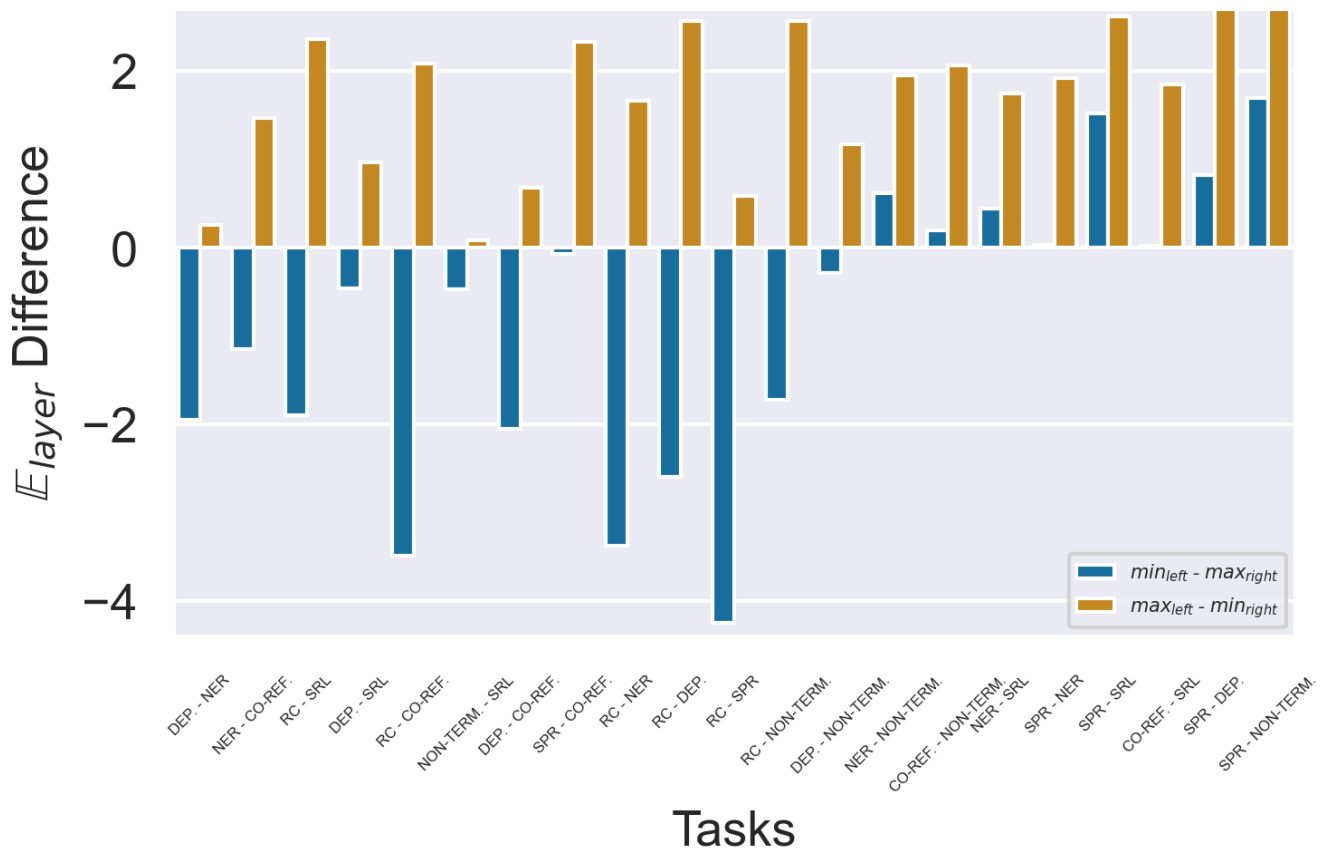


Figure 10: Difference between the minimal possible expected layer of the left task and the maximal possible expected layer of the right task (blue - see legend), and vice-versa (yellow - see legend), for every task-pair.