

# Can You Answer This? - Exploring Zero-Shot QA Generalization Capabilities in Large Language Models

Saptarshi Sengupta <sup>\*1</sup>, Shreya Ghosh<sup>1</sup>, Preslav Nakov<sup>2</sup>, and Prasenjit Mitra<sup>1</sup>

<sup>1</sup>College of IST, The Pennsylvania State University, USA

<sup>2</sup>Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE

## Abstract

The buzz around Transformer-based language models (TLM) is well founded owing to their impressive performance on a number of tasks. However, when applied to areas that require specialized knowledge (closed-domain), such as medical, finance, etc., their performance takes drastic hits, sometimes more than their older recurrent/convolutional counterparts. In this paper, we explore *zero-shot* capabilities of large LMs for extractive QA. Our objective is to examine performance change in the face of *domain drift*, i.e., when the target domain data is vastly different in terms of semantic and statistical properties from the source domain, and attempt to explain the subsequent behavior. To this end, **we present two studies**, and plan further experiments later down the road. Our findings indicate certain flaws in the current generation of TLM preventing them from performing well on closed-domain tasks.

## 1 Introduction

The optimism surrounding the achievements of TLM such as BERT [3], RoBERTa [5], and XLNet [12] starts getting tested when we probe their generalization capabilities, i.e. asking whether they are capable of performing well on any domain. According to the established pre-training + fine-tuning framework, these models should be fine-tuned on every new dataset they are to be applied to in order to achieve SOTA performance. However, an increasing number of studies [6, 8, 4] have started questioning the limits of this approach and exploring *no-fine-tune* or *zero-shot* and *few-shot* transfer learning settings.

Closed-domain datasets (CDD) pose challenges that TLM are not equipped to tackle. First, the

language used to describe phenomena in their space is quite dense i.e. filled with jargon that typically do not appear in general purpose corpora (out-of-vocabulary issue). Second, the number of samples in these datasets aren't nearly enough to provide competitive performance via fine-tuning since to work well, TLM need numerous samples to learn a decent enough approximation of the data. Third, the statistical demands of these datasets such as answer and context lengths are much more than what these TLM are used to seeing.

The above properties, while useful for understanding the nuances of CDD, do not completely explain *why* there exists performance discrepancies due to domain shift. For example, even if we use models specifically tailored for a given domain, say biomedical, we see that their performance is still quite poor; ex. SciBERT [1] when applied to COVID-QA [7] gives mediocre zero-shot performance  $\approx 0.45$  F1 & 0.25 EM, and even when fine-tuned does not improve significantly, 0.54 F1 & 0.29 EM [10].

Thus, we investigate two hypothesis (in this abstract), based on a study of related datasets and literature, to understand the root cause for such performance discrepancy. They are **generated answer length analysis** and **semantic similarity between different senses of polysemous domain terms**. We have observed that CDD usually demand longer answer spans to be generated as opposed to simple factoid based QA. This makes sense seeing as we need detailed answers to questions such as *Why is remdesivir ineffective for COVID-19?* as opposed to simple factoid based questions. As such, we wanted to see whether these models are capable of producing longer spans. The reasoning behind the second test is, typically, it is expected that such datasets would showcase only one sense of a polysemous word i.e. the sense which is *usually* talked about in that space. Thus, we needed to see if these models are capable

---

\*Corresponding Author: sks6765@psu.edu

Word Model \ Sense	Server			Java			Windows			following			Min	Max	Avg.
	host (s1)	waiter (s2)	(s1, s2)	software (s1)	coffee (s2)	(s1, s2)	software (s1)	framework (s2)	(s1, s2)	reference (s1)	pursue (s2)	(s1, s2)			
BERT	0.78	0.8	0.48	0.84	0.67	0.48	0.77	0.71	0.33	0.61	0.58	0.31	0.31	0.84	0.61
RoBERTa	0.94	0.95	0.89	0.94	0.91	0.84	0.93	0.94	0.79	0.91	0.92	0.84	0.79	0.95	0.90
SciBERT	0.79	0.82	0.71	0.79	0.65	0.59	0.71	0.85	0.62	0.67	0.69	0.55	0.55	0.85	0.70
SenseBERT	0.88	0.89	0.75	0.90	0.69	0.53	0.83	0.92	0.76	0.72	0.81	0.52	0.52	0.92	0.77

Table 1: Avg. Semantic Similarity between contextualized (s)senses of the domain terms as found in TechQA.

Model	SQuADv1				Model	TechQA			
	Gold	Predicted	EM	F1		Gold	Predicted	EM	F1
BIDAF	18.73	<b>25.31</b>	65.73	75.98	BIDAF	156.79	<b>4302.93</b>	<b>32.23</b>	<b>39.45</b>
QANet		23.74	26.3	36.81	QANet		387.2	3.96	7.65
BERT		18.18	80.95	88.25	BERT		18.42	1.61	6.35
RoBERTa		18.03	<b>82.73</b>	<b>90.04</b>	RoBERTa		26.89	1.94	4.68

Table 2: Average no. of characters in the answer spans generated for questions in the validation set of SQuAD. **Note, zero-shot performance (EM & F1) is also provided alongside.**

Table 3: Average no. of characters in the answer spans generated for questions in the validation set of TechQA.

of distinguishing between the desired and auxiliary senses of a word. An inability to do so would indicate a deep semantic drawback.

## 2 Answer Length Analysis

For this experiment, we computed the average number of characters present in the generated answer spans on the validation set of TechQA [2] (our CDD) and SQuAD [9] (our ODD). SQuAD was chosen as the ODD since there exist many fine-tuned models, eliminating the need for retraining. We count characters instead of tokens as we wanted a scheme independent of tokenization strategy (ex. BPE, wordpiece, etc.) We present our results in tables 2 and 3.

From the presented scores, we see that CDD do require longer answer spans of text which the TLM aren’t able to generate. To see whether this is a problem with neural architectures in general, we examine scores from two other variants viz. recurrent (BIDAF [11]) and convolution (QANet [13]). While both are capable of generating incredibly longer spans of text, interestingly, BIDAF performs multitudes better than either TLM on **this CDD** (remains to be seen for other CDD). Further analysis is needed to determine whether this behavior exists across other CDD and/or whether it can achieve high scores simply due to its longer spans because the current F1 metric is based on token overlap. We would be circumspect to linearly correlate performance with predicted sequence length.

## 3 Semantic Similarity Trials

**We conjecture** that polysemous words in CDD only manifest the *dominant* sense in their space. Ex. Given the number of samples, it would be difficult to expect a technical QA dataset to have (m)any instances of the *coffee* sense for *Java* (although this may not always be the case as in *cold* for temperature & condition for biomedical datasets). To test this hypothesis, we create a dataset of polysemous domain terms, appearing in the vocabulary of the TLM and the corpus (TechQA), and associated contexts. As expected, the corpus shows only a single sense of a word & as such, we scraped an online resource, [vocabulary.com](http://vocabulary.com), for contexts for the other sense of the words. In total, we had ten contexts per sense of a given word. We compute average cosine similarity b/w contextualized embeddings of the target word from same & different sense groups. The logic here is that intra/same-group similarity should be high while inter/different-group should be low. According to table 1 BERT *is able* to respect this logic whereas the others, including SenseBERT (surprisingly) is not.

## 4 Conclusion

Acknowledging these results as preliminary, we feel that scores from the presented trials do provide some insight into the poor performance of TLM on CDD. However, we need to run these experiments & more on other datasets to make our claims ironclad.

## References

- [1] I. Beltagy, K. Lo, and A. Cohan. SciBERT: A pretrained language model for scientific text. 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 3615–3620, Hong Kong, China, Nov. 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1371. URL <https://aclanthology.org/D19-1371>.
- [2] V. Castelli, R. Chakravarti, S. Dana, A. Ferritto, R. Florian, M. Franz, D. Garg, D. Khandelwal, J. S. McCarley, M. McCawley, et al. The techqa dataset. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1269–1278, 2020.
- [3] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. 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, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [4] A. Kumar and V. H. C. Albuquerque. Sentiment analysis using xlm-r transformer and zero-shot transfer learning on resource-poor indian language. *Transactions on Asian and Low-Resource Language Information Processing*, 20(5):1–13, 2021.
- [5] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [6] Q. Lyu, H. Zhang, E. Sulem, and D. Roth. Zero-shot event extraction via transfer learning: Challenges and insights. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 322–332, Online, Aug. 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.42. URL <https://aclanthology.org/2021.acl-short.42>.
- [7] T. Möller, A. Reina, R. Jayakumar, and M. Pietsch. COVID-QA: A question answering dataset for COVID-19. In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*, Online, July 2020. Association for Computational Linguistics. URL <https://aclanthology.org/2020.nlpccovid19-acl.18>.
- [8] M. Moradi, K. Blagec, F. Haberl, and M. Samwald. Gpt-3 models are poor few-shot learners in the biomedical domain. *arXiv preprint arXiv:2109.02555*, 2021.
- [9] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, 2016.
- [10] S. Sengupta, C. Heaton, S. Sarkar, and P. Mitra. Leveraging external knowledge resources to enable domain-specific comprehension. Accepted to CoLLAs Workshop, 2022.
- [11] M. Seo, A. Kembhavi, A. Farhadi, and H. Hajishirzi. Bidirectional attention flow for machine comprehension. *arXiv preprint arXiv:1611.01603*, 2016.
- [12] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2019.
- [13] A. W. Yu, D. Dohan, M.-T. Luong, R. Zhao, K. Chen, M. Norouzi, and Q. V. Le. Qanet: Combining local convolution with global self-attention for reading comprehension. *arXiv preprint arXiv:1804.09541*, 2018.