

# QASE Enhanced PLMs: Improved Control in Text Generation for MRC

Anonymous ACL submission

## Abstract

To address the challenges of out-of-control generation in generative models for machine reading comprehension (MRC), we introduce the **Question-Attended Span Extraction (QASE)** module. Integrated during the fine-tuning of pre-trained generative language models (PLMs), *QASE* enables these PLMs to match SOTA extractive methods and outperform leading LLMs like GPT-4 in MRC tasks, without significant increases in computational costs.<sup>1</sup>

## 1 Introduction

Machine Reading Comprehension (MRC) is a critical NLP challenge. Mainstream approaches to MRC extract a relevant piece of text from the context in response to a question (Wang et al., 2018; Yan et al., 2019; Chen et al., 2020), but in real-world application, the correct answers often span multiple passages or are implicit (Li et al., 2021). Exploring generative models, in addition to extractive methods, is essential.

Generative models, however, underperform in MRC due to out-of-control generation (Li et al., 2021). This leads to two main challenges: (1) ill-formed generated answers, containing incomplete or redundant phrases, and (2) factual inconsistency in the generated answers deviating from the correct response. In this paper, we address these by introducing a lightweight **Question-Attended Span Extraction (QASE)** module. We fine-tune multiple open-source generative pre-trained language models (PLMs) on various MRC datasets to assess the module’s efficacy in guiding answer generation. Our contributions include: (1) Developing *QASE* to improve fine-tuned generative PLMs’ quality and factual consistency on MRC tasks, matching SOTA extractive methods and surpassing GPT-4; (2) *QASE* boosts performance without signif-

icantly increasing computational costs, benefiting researchers with limited resources.

## 2 Related Work

Most **current studies on MRC** involve predicting the start and end positions of the answer spans from a given context (Ohsugi et al., 2019; Lan et al., 2019; Bachina et al., 2021; Chen et al., 2022) using encoder-only PLM models such as BERT and XLM-Roberta. To handle the multi-span setting, some studies frame the problem as a sequence tagging task (Segal et al., 2020), and others explore ways to combine models with different tasks (Hu et al., 2019; Lee et al., 2023; Zhang et al., 2023). While these extractive-based methods mainly utilize encoder-only models, there is also research focuses on using generative language models (Yang et al., 2020; Li et al., 2021; Su et al., 2022).

**Retrieval-augmented text generation (RAG)** augments the input of PLMs with in-domain (Gu et al., 2018; Weston et al., 2018; Saha and Srihari, 2023) or external knowledge (Su et al., 2021; Xiao et al., 2021) to control the quality and factual consistency of generated content. It has become a new text generation paradigm in many NLP tasks (Li et al., 2022b), such as dialogue response generation (Wu et al., 2021; Liu et al., 2023b) and machine translation (He et al., 2021; Zhu et al., 2023). However, not much work focuses on selective MRC. Our approach diverges from RAG as it directly fine-tunes the weights of the PLMs rather than altering the input to the PLMs with additional information.

## 3 Method

**Question-Attended Span Extraction** To guide text generation, we use *QASE*, a question-attended span extraction module, during fine-tuning the generative PLMs. *QASE* focuses model attention on potential answer spans within the original context. We cast span extraction as a sequence tagging prob-

<sup>1</sup>Our code is available at [this anonymous repo link](#).

lem and employ the Inside-Outside (IO) tagging schema, where each sequence token is tagged as ‘inside’ (*I*) if part of a relevant span, or ‘outside’ (*O*) if not. This schema works well for both single- and multi-span extraction settings, achieving comparable or even better performance than the well-known BIO tagging format (Huang et al., 2015), as shown by Segal et al. (2020).

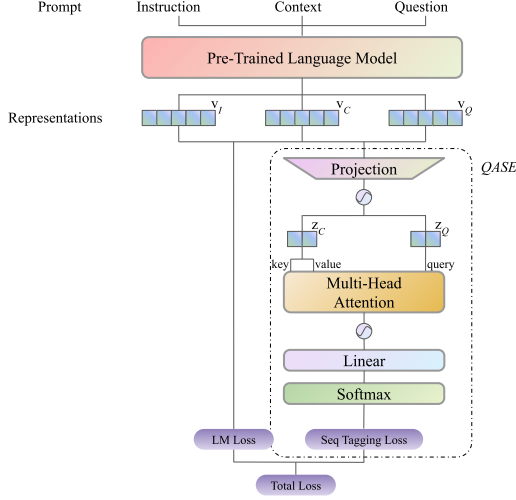


Figure 1: *QASE*-enhanced model architecture

The architecture of our model is shown in Figure 1. An input context and question pair and an instruction are first tokenized and fed into the PLM. The hidden states output from the PLM is then passed through projection layers to produce embeddings  $z_i = \text{ReLU}(W_{\text{proj}}v_i + b_{\text{proj}})$ , where  $v_i \in R^d$  is the PLM output hidden state of the  $i^{\text{th}}$  token.

To learn context tokens representations in relation to specific questions, we employ a **multi-head attention** mechanism (*MHA*). Each head in *MHA* focuses to different aspects of the context as it relates to the question, using question embeddings as the query and context embeddings as key-value pairs. This mechanism aligns the context token representations with the specifics of the queried question. The projected embeddings  $z_i$  are passed through *MHA*, and subsequently channeled through a linear layer and a softmax layer to compute  $p_i = \text{softmax}(W_{\text{lin}} \cdot \text{MHA}(z_i) + b_{\text{lin}})$ , which denotes the probability of the  $i^{\text{th}}$  token being inside the answer spans. We then compute the sequence tagging loss using the cross entropy loss  $L_{QASE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=0}^1 y_{ij} \log(p_{ij})$ , where  $j \in 0, 1$  corresponds to class *O* and class *I*, and  $y_{ij}$  is a binary value indicating whether the  $i^{\text{th}}$  token belongs to class  $j$ .

**Fine-Tuning and Inference** We fine-tune the

PLMs using multi-task learning, simultaneously optimizing both the language modeling loss and sequence tagging loss:  $L = L_{LML} + \beta L_{QASE}$ , where  $\beta$  is a hyper-parameter that controls the weight of the span extraction task. This approach enhances the PLMs’ ability to generate answers well-founded in the context and relevant answer spans. During inference, only the generation component of the fine-tuned model is employed.

## 4 Experiments

**Datasets and Metrics** We utilize these 3 MRC datasets. (1) **SQuAD** (Rajpurkar et al., 2016): A benchmark dataset consisting of 100K+ questions with single-span answers. We use SQuAD v1.1. Since the official evaluation on v1.1 has long been ended, we report our results on the official v1.1 development set. (2) **MultiSpanQA** (Li et al., 2022a): This dataset consists of over 6.5k question-answer pairs. Unlike most existing single-span answer MRC datasets, MultiSpanQA focuses on multi-span answers. (3) **Quoref** (Dasigi et al., 2019): A benchmark dataset containing more than 24K questions, with most answers being single-span and  $\sim 10\%$  being multi-span. Following the conventions of the datasets’ official leaderboards (listed in A.1), we employ exact match (EM) and partial match (Overlap) F1 scores as metrics on MultiSpanQA, and exact match percentage and macro-averaged F1 score on SQuAD and Quoref.

**Experimental Setup** To evaluate the effectiveness of *QASE* independent of any specific language model, we experiment with multiple open-source LLMs. These include both decoder-only LLMs, such as Llama 2 (Touvron et al., 2023) and Alpaca (Taori et al., 2023), and an encoder-decoder model, Flan-T5 (Chung et al., 2022). For Llama 2 and Alpaca, we fine-tune the pre-trained 7B version using LoRA (Hu et al., 2021) and instruction-tuning (see A.4 for instruction templates). For Flan-T5 family models, we fine-tune the small, the base, and the large versions. The trainable parameters for each model is provided in Table 2.

We set the hyper-parameters  $\beta = 1$  and the learning rate  $lr = 1e-4$ . For LoRA fine-tuning applied to Llama 2 and Alpaca models, we specify a rank  $r = 8$ ,  $\alpha = 32$ , and a dropout rate of 0.05. The methodology for selecting these hyper-parameters is detailed in A.2. We train all our models on single GPUs, using a batch size of 2-4 depending on the VRAM of the respective GPUs. We use four types

		Llama2	Alpaca	Flan-T5-Small	Flan-T5-Base	Flan-T5-Large
<b>SQuAD</b> (EM   F1)	no <i>QASE</i>	36.68   47.06	27.88   43.95	77.33   85.51	82.09   89.56	83.16   90.71
	<i>QASE</i>	<b>37.22   47.69</b>	<b>37.31   47.62</b>	<b>77.66   85.90</b>	<b>82.20   90.24</b>	<b>84.13   91.70</b>
<b>MultiSpanQA</b> (EM F1   Overlap F1)	no <i>QASE</i>	50.93   68.14	<b>52.73</b>   69.10	<b>59.13</b>   76.49	64.66   81.41	<b>67.41</b>   83.09
	<i>QASE</i>	<b>51.75   70.39</b>	52.20   <b>70.01</b>	59.08   <b>77.10</b>	<b>64.87   81.50</b>	66.92   <b>84.22</b>
<b>Quoref</b> (EM   F1)	no <i>QASE</i>	45.52   52.09	-	58.21   63.30	72.77   80.90	75.17   80.49
	<i>QASE</i>	<b>54.28   60.44</b>	-	<b>60.70   66.88</b>	<b>75.17   81.18</b>	<b>76.19   82.13</b>

Table 1: Performance of fine-tuned PLMs with or without *QASE* on each dataset.

	Trainable Parameters		
	no <i>QASE</i>	<i>QASE</i>	$\Delta$ params
<b>Llama2/Alpaca with LoRA</b>	4.2M	7.3M	3.1M
<b>Flan-T5-Small</b>	77.0M	78.2M	1.3M
<b>Flan-T5-Base</b>	247.6M	248.9M	1.4M
<b>Flan-T5-Large</b>	783.2M	784.7M	1.5M

Table 2: Trainable parameters of experimented models.

of GPUs: A40, A10, A5500, and A100. Models are trained for 3 epochs or until convergence.

**Experiment Results** To evaluate the efficacy of the *QASE*, we examine the performance of various PLMs fine-tuned with and without *QASE*, as shown in Table 1. Generally, models fine-tuned with *QASE* outperform those fine-tuned without it. In particular, for SQuAD, *QASE*-enhanced model demonstrate an EM percentage increase of up to 33.8% and an F1 score upsurge of up to 8.4% over vanilla fine-tuned models. For MultiSpanQA, there is an improvement of up to 1.6% in the EM F1 and up to 3.3% in the overlap F1. Likewise, on Quoref, there is an improvement of up to 19.2% in the EM percentage and up to 16.0% in the F1 score. These results show that, by employing *QASE*, generative-based PLMs can be fine-tuned to produce well-formed, context-grounded, and better-quality answers in MRC tasks compared to the vanilla fine-tuning approach. For reference, we also compare the fine-tuned PLMs to their corresponding PLMs in zero-shot settings, as presented in Appendix A.3.

**Computational Costs** Table 2 shows that integrating *QASE* slightly raises the number of trainable parameters in PLMs, with the increase dependent on the models’ hidden sizes. Significantly, for the largest model, Flan-T5-Large, *QASE* adds just 0.2% more parameters, indicating that *QASE* enhances the capabilities of fine-tuned PLMs in MRC without major increase in computational resources.

**Model Comparisons** Our top model, Flan-T5-Large<sub>*QASE*</sub>, is further benchmarked against leading models on each dataset’s official leaderboard, alongside zero-shot GPT-3.5-Turbo and GPT-4. GPT-3.5-Turbo stands as one of OpenAI’s most

efficient models in terms of capability and cost, while GPT-4 shows superior reasoning abilities (Liu et al., 2023c). Studies indicate their superiority over traditional fine-tuning methods in most logical reasoning benchmarks (Liu et al., 2023a). The prompts used to query the GPT variants are detailed in Appendix A.4. On SQuAD, as showed in Table 3, Flan-T5-Large<sub>*QASE*</sub> surpasses human performance, equaling the NLNet model. Additionally, it surpasses GPT-4 by 113.8% on the exact match score and 32.6% on F1. On MultiSpanQA, Table 4

	EM	F1 $\uparrow$
GPT-3.5-Turbo	36.944	65.637
GPT-4	39.347	69.158
Human Performance	82.304	91.221
BERT-Large (Devlin et al., 2019)	84.328	91.281
MSRA NLNet (ensemble)	<b>85.954</b>	91.677
Flan-T5-Large <sub><i>QASE</i></sub>	84.125	<b>91.701</b>

Table 3: Flan-T5-Large<sub>*QASE*</sub> and baselines on SQuAD.

shows that Flan-T5-Large<sub>*QASE*</sub> outperforms LIQUID (Lee et al., 2023), which currently ranks #1 on the leaderboard, with respect to the overlap F1 score. Moreover, it surpasses GPT-4 by 4.5% on the exact match F1 and 1.5% on the overlap F1. On

	EM F1	Overlap F1 $\uparrow$
GPT-3.5-Turbo	59.766	81.866
GPT-4	64.027	82.731
LIQUID (Lee et al., 2023)	<b>73.130</b>	83.360
Flan-T5-Large <sub><i>QASE</i></sub>	66.918	<b>84.221</b>

Table 4: Performance of Flan-T5-Large<sub>*QASE*</sub> and baselines on MultiSpanQA.

	EM	F1 $\uparrow$
GPT-3.5-Turbo	50.22	59.51
GPT-4	68.07	78.34
CorefRoberta-Large (Ye et al., 2020)	75.80	<b>82.81</b>
Flan-T5-Large <sub><i>QASE</i></sub>	<b>76.19</b>	82.13

Table 5: Performance of Flan-T5-Large<sub>*QASE*</sub> and baselines on Quoref.

Quoref, Table 5 shows that Flan-T5-Large<sub>*QASE*</sub> is comparable to CorefRoberta-Large (Ye et al., 2020), which ranks #9 on the leaderboard, with

a 0.5% higher exact match. Furthermore, it outperforms GPT-4 by 11.9% on the exact match and 4.8% on F1.

All top-performing models on these datasets’ leaderboards, equaling or exceeding Flan-T5-Large<sub>QASE</sub>, are encoder-only extractive models. Therefore, these results demonstrate that *QASE*-enhanced generative PLMs can be fine-tuned to match or exceed the capabilities of SOTA extractive models and outperform leading LLMs in MRC.

**Ablation Studies** To demonstrate the superiority of the *QASE* architecture, we compared Flan-T5-Large<sub>QASE</sub> with vanilla fine-tuned Flan-T5-Large<sub>FT</sub> and Flan-T5-Large<sub>baseline</sub>. The baseline span extraction module lacks the *MHA* component, making it a standard architecture for fine-tuning pre-trained encoders for downstream sequence tagging tasks. We also explored both question-first (*qf*) and context-first prompting strategies, with further details and analysis provided in Appendix A.5, where the model architecture is also illustrated.

Table 6 shows that the baseline-embedded model performs better with a question-first prompting strategy, as Flan-T5-Large<sub>baseline<sub>qf</sub></sub> surpasses Flan-T5-Large<sub>baseline</sub> and Flan-T5-Large<sub>FT<sub>qf</sub></sub>. Conversely, the baseline span extraction module decreases performance in context-first prompting, where Flan-T5-Large<sub>baseline</sub> underperforms compared to Flan-T5-Large<sub>FT</sub>. This suggests that adding an auxiliary span extraction module without careful design can negatively affect instruction fine-tuning. Meanwhile, the *QASE*-enhanced model excels over both vanilla fine-tuned and baseline-embedded models in both prompting scenarios, demonstrating its architectural superiority. Specifically, in context-first setting, Flan-T5-Large<sub>QASE</sub> significantly outperforms Flan-T5-Large<sub>baseline</sub> with a 4.3% higher F1.

	EM	F1 ↑
Flan-T5-Large <sub>baseline</sub>	79.877	87.918
Flan-T5-Large <sub>FT<sub>qf</sub></sub>	80.378	88.176
Flan-T5-Large <sub>baseline<sub>qf</sub></sub>	81.125	89.043
Flan-T5-Large <sub>QASE<sub>qf</sub></sub>	81.485	89.077
Flan-T5-Large <sub>FT</sub>	83.159	90.712
Flan-T5-Large <sub>QASE</sub>	<b>84.125</b>	<b>91.701</b>

Table 6: Performance of vanilla, baseline-, and *QASE*-enhanced fine-tuned Flan-T5-Large on **SQuAD**.

**Factual Consistency** While token-based EM and F1 scores measure the structural quality of generated text, they do not reflect factual accuracy relative to the context. For this we used  $Q^2$  (Hon-

ovich et al., 2021), an automatic metric for assessing factual consistency in generated text, which uses question generation and answering methods over token-based matching. We compared fine-tuned Flan-T5-Large with and without *QASE* in both single-span (SQuAD) and multi-span (MultiSpanQA) answer settings. Table 7 shows that *QASE*-enhanced models consistently outperform the vanilla fine-tuned model. On SQuAD,  $Q^2$  NLI score is improved by 1.0%, and on MultiSpanQA, it is improved by 16.0%. Beyond the  $Q^2$  statistical analysis, our detailed case studies in Appendix A.6 highlight Flan-T5-Large<sub>QASE</sub>’s improved performance. These examples show the model’s better alignment with relevant context, its enhanced understanding of complex sentences, its skill in synthesizing answers from dispersed information, and its superior use of pre-existing real-world knowledge in generating answers.

	Flan-T5-Large	$Q^2$ F1	$Q^2$ NLI
<b>SQuAD</b>	no <i>QASE</i>	42.927	44.983
	<i>QASE</i>	<b>43.624</b>	<b>45.419</b>
<b>MultiSpanQA</b>	no <i>QASE</i>	32.889	31.433
	<i>QASE</i>	<b>34.732</b>	<b>36.452</b>

Table 7:  $Q^2$  scores of fine-tuned Flan-T5-Large with or without *QASE* on each dataset.

## 5 Conclusion and Future Work

In this study, we address out-of-control text generation of generative PLMs in MRC using *QASE*, a lightweight question-attended span extraction module, during the fine-tuning of PLMs. Our experiments show that *QASE*-enhanced PLMs generate better-quality responses with improved formality and factual consistency, matching SOTA extractive models and outperforming GPT-4 by a significant margin on all three MRC datasets. Importantly, *QASE* improves performance without a significant increase in computational costs, benefiting researchers with limited resources.

In the future, we plan to test our model on generative MRC datasets (Nguyen et al., 2016) to further assess its efficacy in more complex scenarios. Another key focus will be evaluating the model’s general ability in answer generation, particularly from the perspective of human perception. This will involve incorporating human annotators in addition to automatic metrics. For a long-term goal, we are looking to expand our work to explore solutions for addressing input- and context-conflicting hallucinations in LLMs.



## Limitations

Due to our limited computational resources, we have been able to perform our experiments on models no larger than Flan-T5-Large. This same constraint led us to only fine-tuning of Llama 2 and Alpaca with LoRA. We note that models based on Llama 2 and Alpaca generally underperform those based on Flan-T5. Apart from the inherent distinctions between decoder-only and encoder-decoder models, and their suitability for different tasks (as seen from the models’ zero-shot performance), a possible factor could be the number of trainable parameters during fine-tuning. Specifically, fine-tuning Llama 2 and Alpaca with LoRA results in only 4.2M trainable parameters, while even the smallest Flan-T5 model provides 77.0M trainable parameters, as shown in Table 2. We acknowledge that many researchers face similar computational resource limitations. Therefore, our research should be very useful, proposing this lightweight module capable of enhancing smaller PLMs to outperform leading LLMs on MRC tasks like these, achieving a balance of effectiveness and affordability.

One foreseeable limitation of our work is the dependency of the fine-tuning process on answer span annotations, since *QASE* works as an auxiliary supervised span extraction module. This reliance on annotated data could potentially limit the model’s broader applicability. A prospective exciting future direction to address this limitation is to develop a semi- or unsupervised module that focuses on selecting relevant spans or rationales within a given context. By integrating this module with our current model, we could significantly improve its generalization capabilities, thereby making it more adaptable and effective across a wider range of scenarios.

One popular method to enhance the formality of answers generated by LLMs is through prompt engineering, paired with few-shot or in-context learning techniques. While these strategies offer great advantages, our ultimate goal is to create a system with broad domain generalization, one that minimizes the need for extensive, calibrated prompt engineering and sample selections for task adaptation. Although developing a robust prompt engineering framework or paradigm is an appealing direction, our current focus diverges from this path. As a long-term goal, we aim for a solution that handles diverse tasks with minimal task-specific tuning.

## References

- Sony Bachina, Spandana Balumuri, and Sowmya Kamath S. 2021. [Ensemble ALBERT and RoBERTa for span prediction in question answering](#). In *Proceedings of the 1st Workshop on Document-grounded Dialogue and Conversational Question Answering (DialDoc 2021)*, pages 63–68, Online. Association for Computational Linguistics.
- Kunlong Chen, Weidi Xu, Xingyi Cheng, Zou Xiaochuan, Yuyu Zhang, Le Song, Taifeng Wang, Yuan Qi, and Wei Chu. 2020. Question directed graph attention network for numerical reasoning over text. *arXiv preprint arXiv:2009.07448*.
- Nuo Chen, Linjun Shou, Ming Gong, and Jian Pei. 2022. From good to best: Two-stage training for cross-lingual machine reading comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10501–10508.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Pradeep Dasigi, Nelson F. Liu, Ana Marasović, Noah A. Smith, and Matt Gardner. 2019. [Quoref: A reading comprehension dataset with questions requiring coreferential reasoning](#). 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 5925–5932, Hong Kong, China. Association for Computational Linguistics.
- 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.
- Jiatao Gu, Yong Wang, Kyunghyun Cho, and Victor OK Li. 2018. Search engine guided neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Qiuxiang He, Guoping Huang, Qu Cui, Li Li, and Lemao Liu. 2021. [Fast and accurate neural machine translation with translation memory](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3170–3180, Online. Association for Computational Linguistics.
- Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. *Q<sup>2</sup>: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering*. *arXiv preprint arXiv:2104.08202*.

- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Minghao Hu, Yuxing Peng, Zhen Huang, and Dongsheng Li. 2019. A multi-type multi-span network for reading comprehension that requires discrete reasoning. 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 1596–1606, Hong Kong, China. Association for Computational Linguistics.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *arxiv 2015. arXiv preprint arXiv:1508.01991*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Seongyun Lee, Hyunjae Kim, and Jaewoo Kang. 2023. Liquid: A framework for list question answering dataset generation. *arXiv preprint arXiv:2302.01691*.
- Chenliang Li, Bin Bi, Ming Yan, Wei Wang, and Songfang Huang. 2021. Addressing semantic drift in generative question answering with auxiliary extraction. 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 942–947, Online. Association for Computational Linguistics.
- Haonan Li, Martin Tomko, Maria Vasardani, and Timothy Baldwin. 2022a. MultiSpanQA: A dataset for multi-span question answering. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1250–1260, Seattle, United States. Association for Computational Linguistics.
- Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022b. A survey on retrieval-augmented text generation. *arXiv preprint arXiv:2202.01110*.
- Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023a. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*.
- Shuai Liu, Hyundong Cho, Marjorie Freedman, Xuezhe Ma, and Jonathan May. 2023b. RECAP: Retrieval-enhanced context-aware prefix encoder for personalized dialogue response generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8404–8419, Toronto, Canada. Association for Computational Linguistics.
- Xiao Liu, Junfeng Yu, Yibo He, Lujun Zhang, Kaiyichen Wei, Hongbo Sun, and Gang Tu. 2023c. System report for CCL23-eval task 9: HUST1037 explore proper prompt strategy for LLM in MRC task. In *Proceedings of the 22nd Chinese National Conference on Computational Linguistics (Volume 3: Evaluations)*, pages 310–319, Harbin, China. Chinese Information Processing Society of China.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human generated machine reading comprehension dataset. *choice*, 2640:660.
- Yasuhito Ohsugi, Itsumi Saito, Kyosuke Nishida, Hisako Asano, and Junji Tomita. 2019. A simple but effective method to incorporate multi-turn context with BERT for conversational machine comprehension. In *Proceedings of the First Workshop on NLP for Conversational AI*, pages 11–17, Florence, Italy. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. 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, Austin, Texas. Association for Computational Linguistics.
- Sougata Saha and Rohini Srihari. 2023. ArgU: A controllable factual argument generator. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8373–8388, Toronto, Canada. Association for Computational Linguistics.
- Elad Segal, Avia Efrat, Mor Shoham, Amir Globerson, and Jonathan Berant. 2020. A simple and effective model for answering multi-span questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3074–3080, Online. Association for Computational Linguistics.
- Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. 2022. Read before generate! faithful long form question answering with machine reading. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 744–756, Dublin, Ireland. Association for Computational Linguistics.
- Yixuan Su, Yan Wang, Deng Cai, Simon Baker, Anna Korhonen, and Nigel Collier. 2021. Prototype-to-style: Dialogue generation with style-aware editing on retrieval memory. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:2152–2161.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca).

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Wei Wang, Ming Yan, and Chen Wu. 2018. [Multi-granularity hierarchical attention fusion networks for reading comprehension and question answering](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1705–1714, Melbourne, Australia. Association for Computational Linguistics.

Jason Weston, Emily Dinan, and Alexander Miller. 2018. [Retrieve and refine: Improved sequence generation models for dialogue](#). In *Proceedings of the 2018 EMNLP Workshop SCAI: The 2nd International Workshop on Search-Oriented Conversational AI*, pages 87–92, Brussels, Belgium. Association for Computational Linguistics.

Zequ Wu, Michel Galley, Chris Brockett, Yizhe Zhang, Xiang Gao, Chris Quirk, Rik Koncel-Kedziorski, Jianfeng Gao, Hannaneh Hajishirzi, Mari Ostendorf, et al. 2021. A controllable model of grounded response generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14085–14093.

Fei Xiao, Liang Pang, Yanyan Lan, Yan Wang, Huawei Shen, and Xueqi Cheng. 2021. [Transductive learning for unsupervised text style transfer](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2510–2521, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ming Yan, Jiangnan Xia, Chen Wu, Bin Bi, Zhongzhou Zhao, Ji Zhang, Luo Si, Rui Wang, Wei Wang, and Haiqing Chen. 2019. A deep cascade model for multi-document reading comprehension. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7354–7361.

Junjie Yang, Zhuosheng Zhang, and Hai Zhao. 2020. Multi-span style extraction for generative reading comprehension. *arXiv preprint arXiv:2009.07382*.

Deming Ye, Yankai Lin, Jiaju Du, Zhenghao Liu, Peng Li, Maosong Sun, and Zhiyuan Liu. 2020. [Coreferential Reasoning Learning for Language Representation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7170–7186, Online. Association for Computational Linguistics.

Chen Zhang, Jiuheng Lin, Xiao Liu, Yuxuan Lai, Yansong Feng, and Dongyan Zhao. 2023. How many answers should i give? an empirical study of multi-answer reading comprehension. *arXiv preprint arXiv:2306.00435*.

Wenhao Zhu, Jingjing Xu, Shujian Huang, Lingpeng Kong, and Jiajun Chen. 2023. [INK: Injecting kNN](#)

[knowledge in nearest neighbor machine translation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15948–15959, Toronto, Canada. Association for Computational Linguistics.

578  
579  
580  
581  
582



## A Detailed Experiment Setup and Results

### A.1 Dataset Leaderboard

Below are the official leaderboards all the datasets we refer to:

<b>SQuAD</b>	<a href="https://rajpurkar.github.io/SQuAD-explorer/">https://rajpurkar.github.io/SQuAD-explorer/</a>
<b>MultiSpanQA</b>	<a href="https://multi-span.github.io/">https://multi-span.github.io/</a>
<b>Quoref</b>	<a href="https://leaderboard.allenai.org/quoref/submissions/public">https://leaderboard.allenai.org/quoref/submissions/public</a>

Table 8: Dataset official leaderboards.

### A.2 Hyper-Parameter Selection

In this section, we outline the process for selecting the hyper-parameter  $\beta$  and detail our approach to LoRA fine-tuning.

For selecting  $\beta$ , we use a grid search method, exploring values from 0.5 to 2 in increments of 0.1, on 30% of the MultiSpanQA training dataset. This process leads to the determination that  $\beta = 1$  empirically yield the best performance, hence it is selected for use in our experiments.

To select the learning rate  $lr$ , we conduct a grid search, testing values from  $\{1e-5, 5e-5, 1e-4, 5e-4, 1e-3\}$  on 30% of the MultiSpanQA training dataset. Empirically, the value  $1e-4$  demonstrates the best performance and is therefore chosen for our experiments. This selection is in agreement with the default  $lr$  value used in Meta’s official Llama 2 fine-tuning recipe<sup>2</sup>.

In the case of LoRA fine-tuning, we follow the established methodology as outlined by Hu et al. (2021). This involves applying LoRA to Llama 2 and the pre-trained Alpaca models by freezing their pre-trained weights and integrating trainable rank decomposition matrices at every layer of their Transformer structures, aimed at reducing the number of trainable parameters to enhance computational efficiency. We implement this using the PEFT package<sup>3</sup>. The fine-tuning hyper-parameters for LoRA are set according to the default settings specified in Meta’s official Llama 2 fine-tuning recipe<sup>4</sup>, which include a rank  $r = 8$ ,  $\alpha = 32$ , and a dropout rate of 0.05.

<sup>2</sup>Link to the fine-tuning configuration of Meta’s official Llama 2 recipe.

<sup>3</sup>Link to the Hugging Face PEFT implementation.

<sup>4</sup>Link to the LoRA hyper-parameter configuration of Meta’s official Llama 2 recipe.

### A.3 Full Experiment Results

In addition to the highlighted results presented in Section 4, we also compare the fine-tuned PLMs to their corresponding base PLMs in zero-shot settings. The results, presented in Table 9, show that fine-tuning with *QASE* improves performance across all datasets. Specifically, on the SQuAD dataset, models using *QASE* perform up to 5.6 times better in exact match and 3.0 times better in F1 score compared to the original models. On the MultiSpanQA dataset, the exact match improves by up to 124.4 times, and F1 score by up to 3.4 times. Similarly, on the Quoref dataset, the exact match improves by up to 38.4 times, and F1 score by up to 11.2 times with *QASE*.

### A.4 Instruction Templates and Model Prompts

Table 10 provides the instruction and prompt templates used for fine-tuning the PLMs and for zero-shot querying of PLMs and GPT variants across both single- and multi-span answer datasets.

### A.5 Ablation Studies Details

Figure 2 depicts the architecture of the model we use for the ablation studies, with a baseline span extraction module. The baseline span extraction module omits the *MHA* component, typifying a standard architecture for fine-tuning pre-trained encoders for downstream sequence tagging tasks. The baseline-embedded Flan-T5-Large models are fine-tuned with the same configurations as Flan-T5-Large<sub>QASE</sub> including learning rate, weight decay, batch size, epoch number, and GPU type.

We experiment with 2 prompting strategies for ablation studies:

- **Context-first prompting:** The default prompting strategy we utilize for fine-tuning PLMs, both with and without *QASE*. In this setting, the prompt is ordered as "<instruction tokens> <context tokens> <question tokens>".
- **Question-first prompting (*qf*):** Following BERT’s standard fine-tuning procedures. In this setting, the prompt is ordered as "<instruction tokens> <question tokens> <SEP> <context tokens>". <SEP> is a special separator token.



	MultiSpanQA		SQuAD		Quoref	
	EM	F1	EM	F1	EM	F1
Llama2	7.354	34.031	13.443	28.931	5.02	28.91
Llama2 <sub>FT</sub>	50.934	68.140	36.679	47.055	45.52	52.09
Llama2 <sub>QASE</sub>	<b>51.748</b>	<b>70.389</b>	<b>37.219</b>	<b>47.686</b>	<b>54.28</b>	<b>60.44</b>
Alpaca	15.201	42.759	18.259	33.871	-	-
Alpaca <sub>FT</sub>	<b>52.730</b>	69.099	27.881	43.950	-	-
Alpaca <sub>QASE</sub>	52.196	<b>70.008</b>	<b>37.313</b>	<b>47.622</b>	-	-
Flan-T5-Small	0.475	22.539	13.878	28.710	1.58	5.96
Flan-T5-Small <sub>FT</sub>	<b>59.128</b>	76.494	77.332	85.513	58.21	63.30
Flan-T5-Small <sub>QASE</sub>	59.080	<b>77.103</b>	<b>77.663</b>	<b>85.901</b>	<b>60.70</b>	<b>66.88</b>
Flan-T5-Base	4.113	37.694	37.596	51.747	27.08	34.38
Flan-T5-Base <sub>FT</sub>	64.659	81.408	82.090	89.558	72.77	80.90
Flan-T5-Base <sub>QASE</sub>	<b>64.874</b>	<b>81.498</b>	<b>82.204</b>	<b>90.240</b>	<b>75.17</b>	<b>81.18</b>
Flan-T5-Large	13.907	51.501	16.149	37.691	15.96	24.10
Flan-T5-Large <sub>FT</sub>	<b>67.408</b>	83.094	83.159	90.712	75.17	80.49
Flan-T5-Large <sub>QASE</sub>	66.918	<b>84.221</b>	<b>84.125</b>	<b>91.701</b>	<b>76.19</b>	<b>82.13</b>

Table 9: Performance of zero-shot PLMs and fined-tuned PLMs with and without *QASE*.

<b>Fine-tuning</b> PLMs	Instruction: Using the provided context, answer the question with exact phrases and avoid explanations. --- Context: {context} --- Question: {question} --- Answer:
<b>Zero-shot</b> prompting PLMs and GPT variants on <b>single-span</b> answer dataset, SQuAD	Instruction: Using the provided context, answer the question with exact phrases and avoid explanations. --- Context: {context} --- Question: {question} --- Answer:
<b>Zero-shot</b> prompting PLMs and GPT variants on <b>multi-span</b> answer datasets, MultiSpanQA and Quoref	Instruction: Using the provided context, answer the question with exact phrases and avoid explanations. Format the response as follows: ["answer1", "answer2", ...]. --- Context: {context} --- Question: {question} --- Answer:

Table 10: Templates for fine-tuning instructions and zero-shot query prompts

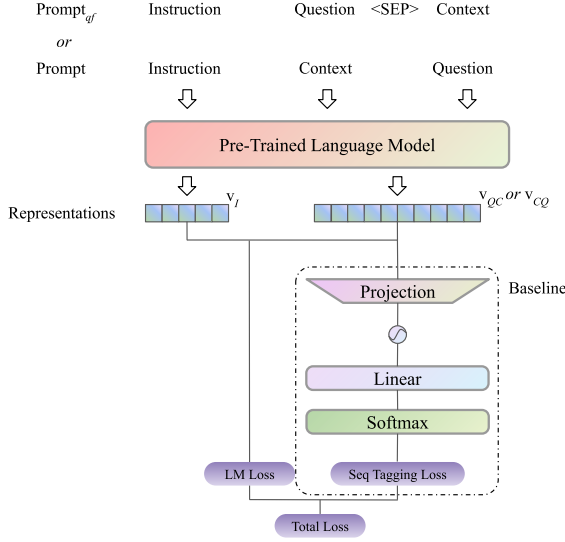


Figure 2: Baseline-embedded model architecture

## A.6 Factual Consistency Case Studies

In Section 4, we demonstrate that the Flan-T5-Large model, when fine-tuned with  $QASE$ , produces answers with greater factual accuracy in relation to the context compared to its counterpart fine-tuned without  $QASE$ . Specifically, we observe a 1.0% improvement in the  $Q^2$  score on the SQuAD dataset and a significant 16.0% increase on Multi-SpanQA. This section includes examples to further illustrate  $QASE$ 's effectiveness.

Table 11 showcases that Flan-T5-Large $_{QASE}$  more accurately identifies the key focus of the question and locates the pertinent factual information within the context, with the aid of the  $QASE$  module. For instance, in **Sample 1**, Flan-T5-Large $_{QASE}$  correctly interprets the question as seeking the age difference between Newton and Manning, rather than the age of either individual, and accordingly provides the accurate answer. In contrast, Flan-T5-Large $_{FT}$  mistakenly provides Newton's age as the answer. Similarly, in **Sample 2**, Flan-T5-Large $_{QASE}$  accurately discerns that the question pertains to Thoreau's claim regarding the majority, generating in the correct answer, whereas Flan-T5-Large $_{FT}$  misguidedly responds with Thoreau's political philosophy.

Flan-T5-Large $_{QASE}$  also shows a notable improvement in comprehending complex, lengthy sentences and synthesizing answers from information that is sparsely distributed across multiple spans requiring logical processing. This capability is particularly valuable when the answer to a question

### Sample 1

**Context:** This was the first Super Bowl to feature a quarterback on both teams who was the #1 pick in their draft classes. Manning was the #1 selection of the 1998 NFL draft, while Newton was picked first in 2011. The matchup also pits the top two picks of the 2011 draft against each other: Newton for Carolina and Von Miller for Denver. Manning and Newton also set the record for the largest **age difference** between opposing Super Bowl quarterbacks at 13 years and 48 days (Manning was 39, Newton was 26).

**Question:** What was the **age difference** between Newton and Manning in Super Bowl 50?

**Gold Answer:** 13 years and 48 days

Flan-T5-Large $_{QASE}$ Generation	13 years and 48 days
Flan-T5-Large $_{FT}$ Generation	26

### Sample 2

**Context:** However, this definition is disputed by Thoreau's political philosophy, which contrasts the conscience with the collective. The individual is the ultimate arbiter of right and wrong. Beyond this, since only individuals act, only they can commit injustices. When the government knocks on the door, it is an individual in the guise of a postman or tax collector whose hand meets the wood. Before Thoreau's imprisonment, when a perplexed tax collector openly pondered how to deal with his refusal to pay, Thoreau had advised, "Resign." If a man chose to be an agent of injustice, then Thoreau insisted on confronting him with the reality that he was making a choice. But if the government is "the voice of the people," as often claimed, shouldn't that voice be heeded? Thoreau acknowledges that the government may represent the will of the majority but it might also merely reflect the desires of elite politicians. Even a good government is "liable to be abused and perverted before the people can act through it." Furthermore, even if a government did express the voice of the people, this fact would not obligate the obedience of individuals who dissent. **The majority may be powerful but it is not necessarily right.** What, then, is the appropriate relationship between the individual and the government?

**Question:** What did Thoreau claim about **the majority**?

**Gold Answer:** not necessarily right

Flan-T5-Large $_{QASE}$ Generation	it is not necessarily right
Flan-T5-Large $_{FT}$ Generation	conscience vs. the collective

Table 11: Comparisons of model attention alignment with question key aspects and relevant factual context between Flan-T5-Large $_{QASE}$  and Flan-T5-Large $_{FT}$ .

does not directly stem from a single phrase. Table 12 provides examples of such instances. In **Sample 3**, the model needs to recognize that ESPN Deportes is the exclusive broadcaster in Spanish and that CBS, although mentioned, does not offer Spanish-language broadcasting. Combining these facts leads to the correct answer, that ESPN Deportes is the network that broadcast the game in Spanish. Flan-T5-Large<sub>QASE</sub> accurately generates this answer, whereas Flan-T5-Large<sub>FT</sub> incorrectly answers with "CBS", likely due to confusion caused by the complex sentence structures and dispersed information. Similarly, in **Sample 4**, Flan-T5-Large<sub>QASE</sub> correctly identifies the question as seeking the name of the force related to a potential field between two locations. It successfully locates the relevant long sentence, deconstructs, and comprehends it to produce the correct answer, in contrast to Flan-T5-Large<sub>FT</sub>, which incorrectly selects the first phrase mentioning "force". In **Sample 5**, the question asks for the class most commonly not ascribed to the graph isomorphism problem. The model needs to deduce from the context that "it is widely believed that the polynomial hierarchy does not collapse to any finite level", implying "graph isomorphism is not NP-complete". Once again, Flan-T5-Large<sub>QASE</sub> arrives at the correct conclusion, while Flan-T5-Large<sub>FT</sub> does not.

While our primary evaluation focuses on the model's proficiency in deriving answers from provided contexts, we also note that *QASE* enhances the model's capacity to leverage real-world knowledge acquired during its pre-training phase. This improvement is attributed to *QASE*'s ability to better align the model's focus on parts of the context that are relevant to the questions asked. Table 13 presents an example of this phenomenon. In **Sample 6**, when asked about the California venue considered for the Super Bowl, Flan-T5-Large<sub>QASE</sub> correctly associates the San Francisco Bay Area with California, thus producing the accurate answer. On the other hand, Flan-T5-Large<sub>FT</sub> erroneously identifies a stadium in Miami as the answer. This example illustrates how *QASE* not only improves context-based answer generation but also the model's application of pre-existing real-world knowledge to the questions posed.

### Sample 3

**Context:** On December 28, 2015, ESPN Deportes announced that they had reached an agreement with CBS and the NFL to be the exclusive Spanish-language broadcaster of the game, marking the third dedicated Spanish-language broadcast of the Super Bowl. Unlike NBC and Fox, CBS does not have a Spanish-language outlet of its own that could broadcast the game (though per league policy, a separate Spanish play-by-play call was carried on CBS's second audio program channel for over-the-air viewers). The game was called by ESPN Deportes' Monday Night Football commentary crew of Alvaro Martin and Raul Allegre, and sideline reporter John Sutcliffe. ESPN Deportes broadcast pre-game and post-game coverage, while Martin, Allegre, and Sutcliffe contributed English-language reports for ESPN's SportsCenter and Mike & Mike.

**Question:** Which network broadcast the game in Spanish?

**Gold Answer:** ESPN Deportes

Flan-T5-Large <sub>QASE</sub> Generation	ESPN Deportes
Flan-T5-Large <sub>FT</sub> Generation	CBS

### Sample 4

**Context:** A conservative force that acts on a closed system has an associated mechanical work that allows energy to convert only between kinetic or potential forms. This means that for a closed system, the net mechanical energy is conserved whenever a conservative force acts on the system. The force, therefore, is related directly to the difference in potential energy between two different locations in space, and can be considered to be an artifact of the potential field in the same way that the direction and amount of a flow of water can be considered to be an artifact of the contour map of the elevation of an area.

**Question:** What is the force called regarding a potential field between two locations?

**Gold Answer:** an artifact

Flan-T5-Large <sub>QASE</sub> Generation	an artifact
Flan-T5-Large <sub>FT</sub> Generation	conservative force

### Sample 5

**Context:** The graph isomorphism problem is the computational problem of determining whether two finite graphs are isomorphic. An important unsolved problem in complexity theory is whether the graph isomorphism problem is in P, NP-complete, or NP-intermediate. The answer is not known, but it is believed that the problem is at least not NP-complete. If graph isomorphism is NP-complete, the polynomial time hierarchy collapses to its second level. Since it is widely believed that the polynomial hierarchy does not collapse to any finite level, it is believed that graph isomorphism is not NP-complete. The best algorithm for this problem, due to Laszlo Babai and Eugene Luks has run time  $2O(\sqrt{n \log(n)})$  for graphs with  $n$  vertices.

**Question:** What class is most commonly not ascribed to the graph isomorphism problem in spite of definitive determination?

**Gold Answer:** NP-complete

Flan-T5-Large <sub>QASE</sub> Generation	NP-complete
Flan-T5-Large <sub>FT</sub> Generation	NP-intermediate

11 Table 12: Comparison of Flan-T5-Large<sub>QASE</sub> and Flan-T5-Large<sub>FT</sub> in understanding complex sentence structures.

Sample 6	
<b>Context:</b> The league eventually narrowed the bids to three sites: New Orleans’ Mercedes-Benz Superdome, Miami’s Sun Life Stadium, and the <b>San Francisco Bay Area’s</b> Levi’s Stadium.	
<b>Question:</b> Which <b>California</b> venue was one of three considered for Super Bowl 50?	
<b>Gold Answer:</b> <b>San Francisco Bay Area’s Levi’s Stadium</b>	
Flan-T5-Large <sub>QASE</sub> Generation	San Francisco Bay Area’s Levi’s Stadium
Flan-T5-Large <sub>FT</sub> Generation	Sun Life Stadium

Table 13: Comparison of Flan-T5-Large<sub>QASE</sub> and Flan-T5-Large<sub>FT</sub> in utilizing real-world knowledge.

## B Extended Discussion on Model Performance

In this section, we engage in a detailed discussion on the performance of the Flan-T5 family of models and Llama 2 in MRC tasks. Our aim is to gain insights into the reasons behind the modest zero-shot performance of these large PLMs on MRC tasks, despite their adeptness at handling other complex NLP tasks such as dialogue generation and summarization. Although a comprehensive analysis falls outside the scope of our current study, exploring these performance nuances can provide valuable perspectives on how to potentially enhance the effectiveness of these PLMs on similar tasks.

### B.1 Discussion on Flan-T5 Zero-Shot Performance

We observe that the zero-shot performance of Flan-T5 models across all datasets, including SQuAD, remains low as shown in Table 9, despite being instruct-tuned on the SQuAD dataset during the pre-training phase. This underperformance might stem from the fact that Flan-T5 models, although trained on the <SQuAD, Extractive QA> task, are also trained on a broad spectrum of 1,836 tasks, predominantly focusing on free-form generation, QA, and reasoning tasks (Chung et al., 2022). Consequently, these models are not finely optimized for extractive QA tasks like MRC, especially under metrics like exact match and F1, particularly for the smaller to larger variants under study. The larger XL and XXL variants may exhibit better performance in these tasks. Furthermore, as discussed in the previous sections, generative models, including Llama 2, Alpaca, and GPT variants, generally show limited effectiveness in MRC tasks in

zero-shot settings, underscored by their poorer performance despite having significantly larger model parameters compared to the Flan-T5 variants we experiment with.

To ensure that our zero-shot experiment’s prompts do not adversely affect Flan-T5’s performance, we compare our prompt template, detailed in Table 10, with those Google released for Flan-T5’s instruct-tuning on the SQuAD v1 dataset<sup>5</sup>. Our template, similar to Google’s, differs mainly by including "with exact phrases and avoid explanations." This difference could potentially affect performance, yet our subsequent experiments demonstrate otherwise.

We conduct a series of experiments to assess the zero-shot performance of Flan-T5-Large on SQuAD, using Google released templates for Flan-T5 instruct-tuning. We select three templates of varying complexities, as listed in Table 14. Our results, detailed in Table 14, reveal that our template achieves the highest F1 score. This indicates the lower performance of zero-shot Flan-T5 on SQuAD and similar MRC datasets is expected, even with the original instruct-tuning templates. It supports our hypothesis that, although Flan-T5 is instruct-tuned on SQuAD, its primary strengths are in broader generative question answering and reasoning, rather than specific extractive QA tasks such as MRC, particularly when evaluated by exact match and F1 metrics.

Prompt Template	SQuAD Performance	
	EM	F1
Article: {context} Question: {question} Answer:	7.001	21.717
Answer a question about this article. Article: {context} Question: {question} Answer:	15.875	33.375
Here is a question about this article: Article: {context} What is the answer to this question: Question: {question} Answer:	<b>16.764</b>	35.304
Our Template See Table 10	16.149	<b>37.691</b>

Table 14: Flan-T5-Large zero-shot performance on SQuAD with different prompt templates.

<sup>5</sup>[Link to Flan-T5 instruct-tuning prompt templates.](#)



## B.2 Discussion on Llama 2 Performance

We observe that models based on Llama 2 and Alpaca generally underperform compared to those based on Flan-T5, in both zero-shot and fine-tuned scenarios, with or without *QASE*. This section delves into a detailed discussion of the potential reasons behind this trend.

Firstly, the discrepancy in performance may stem from the inherent structural differences between decoder-only models (Llama 2 and Alpaca) and encoder-decoder models (Flan-T5). Encoder-decoder models are better equipped for tasks that require extensive input processing, such as MRC, making them more apt for these tasks than decoder-only models, which are typically more suited to open-ended QA scenarios. This fundamental distinction partially accounts for Flan-T5’s superior performance in context-based question answering across both zero-shot and fine-tuned settings.

Additionally, the difference in the number of trainable parameters during fine-tuning might contribute to the observed performance gap. Table 2 indicates that fine-tuning Llama 2 and Alpaca with LoRA leads to a significantly lower count of trainable parameters (4.2M) compared to even the smallest Flan-T5 model (77.0M). This disparity in trainable parameters is a crucial factor in explaining why fine-tuned Flan-T5 models, irrespective of the use of *QASE*, outperform Llama 2 and Alpaca models.

While we address these factors, conducting a comprehensive comparison and analysis of different generative model architectures in MRC tasks exceeds the scope of our current study. Nonetheless, we acknowledge that additional factors, such as the specific instruct-fine-tuning of Flan-T5 models on MRC datasets like SQuAD, might also play a role in their enhanced performance over Llama 2 and Alpaca.