

Query Refinement Prompts for Closed-Book Long-Form QA

Reinald Kim Amplayo Kellie Webster Michael Collins
Dipanjan Das Shashi Narayan
Google DeepMind

{reinald, websterk, mjcollins, dipanjand, shashinarayan}@google.com

Abstract

Large language models (LLMs) have been shown to perform well in answering questions and in producing long-form texts, both in few-shot closed-book settings. While the former can be validated using well-known evaluation metrics, the latter is difficult to evaluate. We resolve the difficulties to evaluate long-form output by doing both tasks at once – to do question answering that requires long-form answers. Such questions tend to be *multifaceted*, i.e., they may have ambiguities and/or require information from multiple sources. To this end, we define *query refinement* prompts that encourage LLMs to explicitly express the multifacetedness in questions and generate long-form answers covering multiple facets of the question. Our experiments on two long-form question answering datasets, ASQA and AQuAMuSe, show that using our prompts allows us to outperform fully finetuned models in the closed book setting, as well as achieve results comparable to retrieve-then-generate open-book models.

1 Introduction

Interest in large language models (LLMs) has exploded in recent years due to their wide range of abilities improving the state-of-the-art with just a few or no examples (Brown et al., 2020; Wei et al., 2022a; Chowdhery et al., 2022; Wei et al., 2022b). One task that benefited greatly is closed-book question answering (Roberts et al., 2020) in a few-shot setting, i.e., to produce correct answers to questions without access to passages to read and find the answer. Most of the impressive results, however, are limited to generating short answers, and while previous work has utilized LLMs to generate long-form texts (Elkins and Chun, 2020; Wei et al., 2022b), most of these outputs are difficult to evaluate since they are subjective in nature (i.e., they may have multiple correct and distinctive answers). In this paper, we attempt to evaluate the long-form generation ability of LLMs through long-form question answering (Fan et al., 2019) in the closed-book setting.

Unlike question answering with short answers, questions with long-form answers are naturally multifaceted,

covering many aspects that are required to fully answer the questions (see Table 1 for examples). This introduces several challenges that LLMs would need to resolve:

1. **Ambiguity:** Many questions in a real-world setting are ambiguous in some way; Min et al. (2020) estimated that more than half of Natural Questions (Kwiatkowski et al., 2019) have multiple plausible interpretations. In the first example of Table 1, the question is ambiguous owing to there being multiple types of cars. A system producing a long-form answer to this question should present the multiple answers in a coherent natural language text.
2. **Information Consolidation:** Some questions require consolidating information from multiple sources to be able to fully answer them (Kulkarni et al., 2020). In the second example of Table 1, information regarding the British and Dutch origins of the term may be found in two different documents and LLMs should therefore be able to synthesize knowledge from multiple sources in a coherent manner.
3. **Correctness:** Having a model to answer questions in a closed-book setting requires it to leverage content it has learned during pretraining. While previous work (Roberts et al., 2020; Brown et al., 2020) has shown that LLMs can do so to produce factoid answers, generating long-form answers requires more and complex information to be leveraged.

Our study is based on two long-form question answering benchmarks, ASQA (Stelmakh et al., 2022) and AQuAMuSe (Kulkarni et al., 2020), which focus on queries that may be ambiguous and/or require information from multiple sources. Using these two datasets, our work is the first to show that LLMs are capable of generating long-form answers to complex questions of various types in a closed-book setting. We devise novel *query refinement* prompts that encourage LLMs to express multiple facets of input question and generate multifaceted answers discussing all identified interpretations coherently. Specifically, we identify several different types of multifacetedness in questions and produce a labeled set of query refinement prompts for question answer pairs with a balanced coverage over different

| ASQA Question: when was the first car made in america? | | |
|--|--|--------------------|
| Multiple Facets | When was the first carriage-sized automobile made in America? | 1871 |
| | When was the first running, gasoline-powered car made in America? | 1893 |
| | ... | ... |
| Long-form Answer | <i>The first carriage-sized automobile that could be driven on wagon roads in the United States was steam-powered and invented in 1871 by Dr. J.W. Carhart in Racine, Wisconsin. The first running, gasoline-powered car that was made in America, the Duryea Motor Wagon, was built in 1893. The Studebaker Automobile Company, which started building cars in 1897, sold electric vehicles in 1902, and gasoline vehicles in 1904.</i> | |
| AQuAMuSe Question: where did the term shooting brake come from? | | |
| Multiple Facets | When did the term originate? | early 19th century |
| | As what term did the term originate? | as a British term |
| | ... | ... |
| Long-form Answer | <i>“Shooting-brake” originated as an early 19th century British term for a vehicle used to carry shooting parties with their equipment and game. The etymology of the term brake is uncertain; initially a chassis used to break in horses, and subsequently used to describe a motorized vehicle. It is also possible, that the word ‘brake’ has its origins in the Dutch word ‘brik’ which means ‘cart’ or ‘carriage’.</i> | |

Table 1: Example multifaceted questions from ASQA and AQuAMuSe and their corresponding answers. A system for closed-book question answering needs to understand that some questions have multiple valid answers and synthesise these into a coherent natural language text output.

types. We then introduce an intermediate query refinement step, a generation subtask akin to explanations and reasoning chains (Wei et al., 2022b; Zhou et al., 2022), where the goal is to identify multiple facets of a given question.

We evaluate our long-form answers using ROUGE (measuring stylistic similarity to gold answers; Lin, 2004) and reading comprehension models (measuring correctness; Stelmakh et al., 2022). Using our query refinement prompts in the few-shot prompting and prompt tuning (Lester et al., 2021) settings, we are able to achieve significantly better performance to fully finetuned closed-book systems, as well as comparable performance to open-book retrieve-then-generate systems, both finetuned on the full training dataset. We further show in our analyses that there are certain biases in the annotation of gold-standard answers, such as summary formatting, that few-shot and prompt tuning methods are not able to capture, thus weakens their automatic metric scores. Our human evaluation shows that answers from prompt tuned LLM are preferred more than those of T5 by humans.

2 Related Work

2.1 Prompting in Large Language Models

LLMs have the revolutionary ability to generalize to tasks presented with natural language prompts (Brown et al., 2020; Chowdhery et al., 2022). This ability has been mostly attributed to the large scale of these models and the learning objective of predicting the next token (Brown et al., 2020). The prompting ability of LLMs has been used successfully in text classification, reading comprehension and open-domain question answering tasks. It has also been shown that LLMs improve on complex reasoning tasks by generating intermediate long-form texts in the form of explanations or reason-

ing steps (Wei et al., 2022b; Zhou et al., 2022). However, prompting is still ineffective when tasked to generate long-form outputs as an *end task*, e.g., generating long summaries for summarization (Brown et al., 2020; Chowdhery et al., 2022). Our work is the first to show that LLMs can do long-form text generation through question answering with the help of a refinement step in the prompt.

Our work is related to and inspired by work on reasoning chains in LLMs (Wei et al., 2022b; Zhou et al., 2022; Kojima et al., 2022; Snell et al., 2022), where the goal is to explicitly generate a reasoning or an explanation before producing an answer. Most of these papers focus on arithmetic and commonsense reasoning questions, where reasoning and explanations are obvious. In this paper, we show that such intermediate explanation generation can also be helpful on tasks that implicitly involve multiple steps, such as long-form question answering where question refinement is necessary. Moreover, we are the first to explore structured explanations in the form of a list of answer facets, which is shown in our experiments to be more effective than natural language explanations.

Prompting is just one way of using the LLMs. There are several work (Sanh et al., 2021; Wei et al., 2021; Chowdhery et al., 2022) that attempted to finetune LLMs entirely for text generation tasks, which can be very expensive. Prompt tuning (Lester et al., 2021) is a popular alternative where *soft* prompts are prepended into the input and are finetuned. There are several other alternatives to prompting that show promising results for generating long-form texts, such as prefix tuning (Li and Liang, 2021), adapters (Bapna and Firat, 2019), and several parameter-efficient finetuning techniques (Clive et al., 2021; He et al., 2022; Liu et al., 2022) that introduce new parameters to the model that is updated during training while leaving the LLM parameter fixed. We

show that applying our query refinement step in prompt tuning also improves the performance.

2.2 Long-form Question Answering

Question answering has emerged as a key way to discover and demonstrate advances in large language models, which are showing their skill on increasingly difficult formulations of the task. SQuAD (Rajpurkar et al., 2016) proposed the first large-scale, human-created reading comprehension task and was used to show the promise of neural architectures, which quickly attained human-like performance on the dataset. Since, there has been a proliferation of reading comprehension datasets developed which probe for specific capabilities (Joshi et al., 2017; Choi et al., 2018; Reddy et al., 2019; Rodriguez et al., 2019). The Natural Questions (Kwiatkowski et al., 2019) effort provided a large reading comprehension dataset based on real information-seeking queries to the Google search engine, and has served most recently a basis for the exploration of questions where a simple short answer is not sufficient to address the information need of a complex question.

One response strategy to such questions is a long-form answer, studied here. Both ASQA (Stelmakh et al., 2022) and AQuaMuSe (Kulkarni et al., 2020) require that systems consolidate information from multiple sources to generate multifaceted long-form answers to questions from the Natural Questions. ASQA focuses on the subset of questions labeled in AmbigQA (Min et al., 2020) for which it is possible to enumerate a collection of refinements and factoid answers that should be covered in a long-form answer. On the other hand, AQuaMuSe focuses on questions without short factoid answers, that typically have a looser relationship to one another. We study both so as to understand what prompting strategies work for the different style of reasoning required to do well on each. ELI5 (Fan et al., 2019) is another long-form question answering dataset that was automatically gathered from Reddit threads, but subsequent work (Krishna et al., 2021) has shown problems in its evaluation, including training/validation overlap and gameable metrics.

Finally, our work is also related to query-focused summarization (Dang, 2005; Zhong et al., 2021; Kulkarni et al., 2020), where a set of relevant passages is assumed to be available. AQuaMuSe was developed for this task but, where our experiments are in the closed-book setting, we discard the given passages.

3 Closed-book Long-form Question Answering

3.1 Few-Shot Prompting Formulation

Given a question q , the goal of closed-book long-form question answering is to produce a passage-length text a without access to external context or knowledge (beyond what was seen in pretraining). For the related task of closed-book *factoid* question answering (Rajpurkar

et al., 2016; Joshi et al., 2017), this can be achieved with large language models using a *few-shot prompting* setup (Brown et al., 2020). That is, given k in-context exemplars of question-answer pairs $[(q'_1, a'_1), \dots, (q'_k, a'_k)]$, usually preceded by an instruction, an LLM will output an answer a for question q from knowledge stored in its parameters (Roberts et al., 2020).

When the answer is instead long-form (see examples in Table 1), there are three subtasks that the model need to do to produce an answer: (1) Determining multiple facets of the question (*Facet Identification*), (2) Retrieving multiple answers to the multiple facets of the question (*Multifaceted Question Answering*), and (3) Realizing a long-form text that includes the multiple answers in a coherent manner (*Surface Realization with Information Consolidation*).

In the next sections, we extend the standard few-shot prompt in three different ways to help LLMs explicitly do these steps before arriving to an answer. Firstly, we identify several types of multifacetedness in questions and produce a labeled and balanced set of exemplars (Section 3.2). Next, we introduce a query refinement step in few-shot prompting that instructs the model to explicitly do the intermediate subtasks (Section 3.3). Finally, we dynamically select exemplars to form a k -shot prompt based on similarity (Section 3.4).

3.2 Types of Multifaceted Questions

There are multiple reasons why a long-form answer would be more felicitous than a factoid answer to a question. Table 2 shows six common types of multifaceted questions in the ASQA (Stelmakh et al., 2022) and AQuaMuSe datasets that fall into this category. These are highly related to the ambiguity categories in AmbigQA (Min et al., 2020), which we used as a seed set for exploring the data in this work. To encourage further work in multifaceted question categorization, we detail the criteria we used to determine what type an example demonstrates. We produce a manually labeled set that contains 20 exemplars of each type to form a pool of possible training instances (see Table 2 for examples). When an example exhibited multiple types simultaneously, we opted not to include it as an exemplar.

3.3 Query Refinement Step

Given that long-form question answering is essentially sequentially solving three subtasks, we propose to use an intermediate step that splits the tasks into two parts. This forces the LLM to explicitly resolve the intermediate subtasks before producing a long-form answer. We experimented with three types of refinements:

- **Natural Language Explanations (NL):** A sentence that explains why the question is ambiguous or needs elaboration. This refinement step essentially does *Facet Identification*, i.e., it identifies the multiple facets of a question, which can either be disambiguations of the question, or things that need

| |
|---|
| <p>(a) Conditional: <i>The original question needs to be refined by specifying additional conditions that may be specifications or constraints.</i> <i>Question:</i> When did movies start being made in color? <i>Multifaceted QA Pairs:</i> Q: When was the first film made that utilized any type of color? A: September 1, 1902 Q: When did the first feature length film come out that was made entirely in three-strip Technicolor? A: June 13, 1935 Related AmbigQA categories: Event references, Properties</p> |
| <p>(b) Set-Valued: <i>The answer to the question is a unstructured collection of size two or greater.</i> <i>Question:</i> What are the neighboring countries of South Korea? <i>Multifaceted QA Pairs:</i> Q: What are the neighboring countries to the North of South Korea? A: North Korea Q: What are the neighboring countries to the South of South Korea? A: Japan Related AmbigQA category: \emptyset</p> |
| <p>(c) Time Dependent: <i>The answer depends on the time at which the question was asked, or changed over time in the past.</i> <i>Question:</i> Where was indian independence league formed in 1942? <i>Multifaceted QA Pairs:</i> Q: Where was indian independence league brought together in March 1942? A: Tokyo Q: Where was indian independence league brought together in June 1942? A: Bangkok Conference Related AmbigQA category: Time-dependency</p> |
| <p>(d) Underspecified Reference: <i>There is a noun phrase in the question which may be resolved in multiple ways.</i> <i>Question:</i> When did bat out of hell come out? <i>Multifaceted QA Pairs:</i> Q: When did the album bat out of hell come out? A: October 21, 1977 Q: When did the TV series bat out of hell come out? A: 26 November 1966 Related AmbigQA category: Entity references</p> |
| <p>(e) Underspecified Type: <i>The entity type or sub-type is not specified in the question.</i> <i>Question:</i> Who is the mayor in horton hears a who? <i>Multifaceted QA Pairs:</i> Q: Who plays the mayor in the 2008 film Horton Hears a Who? A: Steve Carell Q: Who is the mayor in the 2008 film Horton Hears a Who? A: Mayor Ned McDodd Related AmbigQA category: Answer types</p> |
| <p>(f) Needs Elaboration: <i>The answer needs to be elaborated to fully answer the question</i> <i>Question:</i> Where did “you can’t have your cake and eat it too” come from? <i>Multifaceted QA Pairs:</i> Q: Where was the early recording of the phrase found? A: in a letter on 14 March 1538 Q: Who sent the letter? A: Thomas, Duke of Norfolk Q: To whom was the letter sent to? A: Thomas Cromwell Q: How was it phrased in the letter? A: “a man cannot have his cake and eat his cake” Related AmbigQA category: \emptyset</p> |

Table 2: Six types of multifacetedness in questions. The first five types are sampled from the ASQA dataset, while the last type is sampled from the AQuAMuSe dataset. For each type, we also include the related AmbigQA (Min et al., 2020) categories if there are any.

to be elaborated to fully answer the question. In previous work, LLMs have been used to generate NL explanations, e.g. for commonsense reasoning (Ji et al., 2020) and jokes (Chowdhery et al., 2022).

- **Lists of QA Pairs (QA):** We use the multifaceted QA pairs as shown in Table 2. This refinement does both *Facet Identification* and *Multifaceted Question Answering*, which means that when producing the answer, the model just needs to consolidate information into a coherent long-form answer. This is related to question answering blueprints (Narayan et al., 2022), which are used as content plans for summarization models, but their usage has not yet been explored in LLMs.
- **Lists of Answer Facets (AF):** Pairings of the multiple facets of the question and their corresponding answer/s in the form “Facet: Answer”. This is

a concise version of the BP refinement, where we remove uninformative tokens such as wh-words and those that are repeated in the original question.

Table 3 shows all three refinement strategies for all identified ambiguity types. The refinement is inserted between the question and the answer in the exemplar. At inference time, the LLM needs to generate both the refinement and the long-form answer. An example prompt with a query refinement step is illustrated in Figure 1 of Appendix A.

3.4 Dynamic Prompting

Finally, we form a k-shot prompt by selecting from our pool of exemplars created in Section 3.2. We do so using dynamic prompting (DP; Rubin et al., 2022), i.e., ranking exemplars $[(\hat{q}_1, \hat{a}_1), \dots, (\hat{q}_e, \hat{a}_e)]$ using the similarity between input question q and candidate exemplar ques-

Question: When did movies start being made in color?

QA Refinement:
 Q: When was the first film made that utilized any type of color? A: September 1, 1902
 Q: When did the first feature length film come out that was made entirely in three-strip Technicolor? A: June 13, 1935
NL Refinement: The answer depends on what is meant by in color (any type of color or three-strip Technicolor).
AF Refinement:
 - any type of color: September 1, 1902
 - three-strip Technicolor: June 13, 1935

Answer: The first film that utilized any type of color was made September 1, 1902. ... The first feature length film that was made entirely in three-strip Technicolor ... came out on June 13, 1935.

Table 3: An example exemplar with QA, NL, and AF refinements. One exemplar in the k -shot prompt would include a question, one of the three refinements, and the answer. At inference time, the model would need to generate both the refinement and the answer. The answer is deliberately shortened for brevity.

tion \hat{q} . Dynamic prompting helps the model generate refinements for q similarly to how refinements are done for a similar question \hat{q} . We use BERTScore (Zhang* et al., 2020) as our similarity metric. The k most similar exemplars are written to the prompt in reverse order, such that the most similar exemplar is written closest to input question q . In our experiments, we primarily used exemplars with questions labeled as ambiguous for ASQA (top 5 in Table 2, for a total of 100 exemplars), and those that are labeled as *Needs Elaboration* for AQuAMuSe (the last type in Table 2, for a total of 20 exemplars). We also experimented combining both ASQA and AQuAMuSe exemplars in Section 5.

4 Experiments on Ambiguous Question Answering

We conducted experiments on two question answering tasks that require long-form answers: Ambiguous Question Answering (ASQA; Stelmakh et al., 2022) and Query-focused Multi-document Summarization (AQuAMuSe; Kulkarni et al., 2020). In this section, we present results on ASQA; the following section presents results on AQuAMuSe.

4.1 Dataset

The ASQA dataset¹ (Stelmakh et al., 2022) is a long-form question answering dataset built on top of the subset of ambiguous questions identified in the AmbigQA dataset (Min et al., 2020), which itself is a subset of the NQ dataset (Kwiatkowski et al., 2019). ASQA consists of 4,353, 948, and 1,015 training, development, and test examples. Each question (e.g., *Who directed Scarface?*) is paired with a list of QA pairs which sig-

¹<https://www.tensorflow.org/datasets/catalog/asqa>

nify disambiguated questions and their corresponding answers (e.g., *Q: Who directed the 1932 film Scarface? A: Howard Hawks, Q: Who directed the 1983 film Scarface? A: Brian de Palma*). Finally, each example also has two human-written long-form answers based on the given the disambiguated QA pairs. Note that the list of disambiguated questions is not given at inference time.

4.2 Evaluation

Following Stelmakh et al. (2022), we compare systems using three metrics:

- **RL** (ROUGE-L; Lin, 2004): Measures the comprehensibility of the system-generated answer with respect to the gold answers. Since there are two gold answers, we take the maximum ROUGE-L. We lowercase system and gold answers, and report on ROUGE-LSum (f-measure) with the default stemmer on.²
- **D-F1** (Disambiguated F1 Accuracy): Measures the correctness of the system-generated answer. Given the gold-standard disambiguated QA pairs, we run a reading comprehension model (RoBERTa; Liu et al., 2019, trained on SQuAD 2.0; Rajpurkar et al., 2018) where the system-generated long-form answer is the context. We then evaluate the number of disambiguated questions that can be answered using the long-form answer as context by calculating the Q1-F1 accuracy.
- **DR** (Disambiguation-ROUGE): The geometric mean of ROUGE-L and D-F1 which penalizes methods that maximize one metric over another.

Finally, we also report the average number of words of the system-generated answers.

4.3 Results

We compared several systems which can be divided into three types: *few-shot prompting*, *prompt tuning*, and *finetuning*. For few-shot prompting, we used five exemplars in the prompt, and compared several prompt configurations using a 540B-sized large language model (PaLM; (Chowdhery et al., 2022)):

- **Random exemplars (rand)**: Select at random five QA pairs from the training dataset and use as exemplars. We report the average results of five runs.
- **Query diversified exemplars (QD)**: Select five QA pairs with questions that have different classes of ambiguity as classified in Section 3.2. We only considered the 100 exemplars with questions labeled as one of five multifacetedness types, which are originally sampled from the ASQA training dataset (i.e., top 5 in Table 2).

²Using pypi package rouge-score.

| Model | #Words | RL | D-F1 | DR |
|------------------------------------|--------|-------------|-------------|-------------|
| <i>Few-shot Prompting</i> | | | | |
| PaLM 540B (rand) | 75.5 | 33.6 | 18.6 | 25.0 |
| PaLM 540B (QD) | 62.1 | 32.3 | 22.5 | 27.0 |
| + NL refinement | 63.6 | 37.1 | 23.9 | 29.8 |
| + QA refinement | 41.9 | 36.0 | 25.3 | 30.2 |
| + AF refinement | 61.9 | 38.0 | 25.5 | 31.1 |
| GPT 175B (rand)* | 21.6 | 14.5 | 7.7 | 10.6 |
| + AF refinement | 46.5 | 33.7 | 17.1 | 24.0 |
| InstructGPT-3 175B (rand) | 40.7 | 35.4 | 25.0 | 29.7 |
| + AF refinement | 39.0 | 34.9 | 23.4 | 28.8 |
| <i>Few-shot Prompt Tuning (PT)</i> | | | | |
| PaLM 540B 100-shot | 62.5 | 39.2 | 25.0 | 31.3 |
| + AF refinement | 53.8 | 40.3 | 25.4 | 32.0 |
| <i>Using the full dataset</i> | | | | |
| T5-Large CB | 62.5 | 33.5 | 7.4 | 15.7 |
| T5-Large OB 1 Passage | 63.0 | 40.3 | 21.2 | 28.2 |
| T5-Large OB 3 Passages | 71.1 | 42.7 | 25.1 | 32.7 |
| T5-Large OB 5 Passages | 71.6 | 43.0 | 26.4 | 33.7 |
| PaLM 540B PT | 64.1 | 40.7 | 27.8 | 33.5 |

Table 4: Evaluation of several systems on the dev set of the ASQA dataset. The best values for each setting are **bold-faced**. We mark systems with an asterisk (*) if they failed to generate answers for at least half the total number of examples. CB and OB refer to closed-book and open-book settings.

- NL/QA/AF refinement: Include a query refinement step and dynamically select exemplars from a pool of exemplars.

The first block of Table 4 reports the results from few-shot prompting. As can be seen, using exemplars with questions of different types of ambiguity significantly improves over the random baseline. All query refinement prompts improve the performance, where the AF refinement performs the best among them.

To check whether our prompts work with other models, we also test our prompts with 175B-sized GPT-3 (Brown et al., 2020) and InstructGPT-3 (Ouyang et al., 2022) models, named `davinci` and `text-davinci-002`, respectively, the latter finetuned further with humans in the loop to better follow user instructions.³ As shown in the second block of Table 4, using standard prompts in GPT-3 fails to generate answers for at least half the total number of examples. Just by adding our proposed configurations, the performance of GPT-3 significantly increases. Interestingly, we do not see the same increasing trend with the InstructGPT-3 model.

Our prompts can also be applied to prompt tuning (Lester et al., 2021) where a set of learned embeddings called *soft prompts* is prepended to the prompt. In par-

³We note that the training data used for `text-davinci-002` is unknown, i.e., we do not know whether the model had access to supervision from long-form question answering datasets during its training.

ticular, we follow the method in Rubin et al. (2022), i.e., we prepend one soft prompt⁴ to the input and fine-tune it using the 100 exemplars we used for dynamic prompting⁵. The third block in Table 4 reports prompt tuning results, where we see a slight improvement when applying our AF refinement prompts.

Finally, we compare our few-shot systems with systems that are trained with the full dataset. In the final block of Table 4, we show results reported in Stelmakh et al. (2022), which are T5-large models in both closed-book (no retrieved passages) and open-book (1/3/5 retrieved passages using Joint Passage Ranker; Min et al., 2021) settings. We also report PaLM 540B prompt tuned using the full dataset. Our best few-shot systems are surprisingly competitive compared to fully finetuned T5 systems, outperforming the closed book system and the open book system with one retrieved passage. Moreover, the correctness of our best systems as measured by the Disambig-F1 score is on par with the open book T5 models. Finally, prompt tuning PaLM using the full dataset performs the best among all systems in terms of correctness, despite having a lower ROUGE-L score than the best T5 system. We believe that open-book models have higher ROUGE-L scores due to the fact that they have access to retrieved passages that they can directly copy, and which may follow the format of the human-generated answers. We discuss this and other annotator biases further in Section 4.5.

4.4 Ablation Studies

We conducted ablation studies on the best few-shot prompting configuration, which is shown in Table 5. In terms of the number of exemplars, while increasing it from 1 to 5 improves the performance, increasing it from 5 to 10 slightly decreases the overall performance. However we see an increase in ROUGE-L and STR-EM, which shows that access to more data increases the ability of the model to copy answer formats. Moreover, increasing the number of dynamic exemplars leads to performance improvements, which is unsurprising. Also, when removing one component (refinement or dynamic prompting), we see a substantial decrease in performance, where the decrease is larger when refinement is not used. Using a different similarity metric for prompt selection does not have significant changes in the model performance, however model-based metrics such as BERTScore (Zhang* et al., 2020) and BLEURT (Sellam et al., 2020) are slightly better than string-based metrics such as BM25 (Robertson et al., 1995). Finally, the increase in performance of using AF query refinement prompts can also be seen when using smaller versions of PaLM.

⁴We tried several soft prompt lengths and found that increasing the length beyond one prompt does not lead to any improvements.

⁵All our prompt tuning experiments use 2000 steps and batch size of 16.

| Model | #Words | RL | D-F1 | DR |
|---|--------|------|------|------|
| <i>Number of exemplars in the prompt</i> | | | | |
| 1-shot | 55.1 | 35.5 | 23.2 | 28.7 |
| 3-shot | 37.9 | 35.6 | 23.7 | 29.1 |
| 5-shot | 61.9 | 38.0 | 25.5 | 31.1 |
| 10-shot | 66.0 | 38.4 | 24.6 | 30.8 |
| <i>Number of dynamic prompt exemplars</i> | | | | |
| 5 | 69.0 | 38.1 | 22.1 | 29.0 |
| 25 | 68.0 | 37.4 | 23.2 | 29.5 |
| 50 | 63.8 | 38.2 | 23.6 | 30.0 |
| 100 | 61.9 | 38.0 | 25.5 | 31.1 |
| <i>Without one component</i> | | | | |
| w/o refinement | 66.3 | 36.2 | 22.9 | 28.8 |
| w/o dynamic prompting | 40.1 | 35.4 | 25.1 | 29.8 |
| <i>Similarity metric for prompt selection</i> | | | | |
| BERTScore | 61.9 | 38.0 | 25.5 | 31.1 |
| BLEURT | 65.3 | 38.4 | 25.1 | 31.0 |
| BM25 | 64.4 | 38.1 | 24.4 | 30.5 |
| <i>Smaller models</i> | | | | |
| PaLM 8B (random) | 84.1 | 23.0 | 9.2 | 14.6 |
| + AF refinement | 63.9 | 32.4 | 10.0 | 18.0 |
| PaLM 62B (random) | 66.6 | 31.3 | 14.5 | 21.3 |
| + AF refinement | 65.3 | 35.9 | 18.0 | 25.4 |

Table 5: Performance of ablated versions of the best few-shot prompting configuration.

4.5 Further Analyses

Question Disambiguation Analysis In the next paragraphs, we discuss annotator biases in ASQA that few-shot systems are not able to capture. The first bias is in the way questions are disambiguated. We observed cases where the majority ambiguity class selected from dynamic prompting was a plausible type of ambiguity for the given question, but the particular disambiguation in ASQA is of another type. This leads to few-shot systems generating an entirely different yet also correct answer to the question. An example of which is shown in the first block of Table 8 in Appendix B. In fact, by just providing how the ASQA dataset expects the question to be disambiguated (prompt template shown in Section A), the performance of our best configuration improves by up to 8.8 DR points (29.6 vs 38.4). Future work on long-form question answering evaluation should explore methods to deal with questions that can be disambiguated in multiple ways.

Summary Format Analysis The second bias is the way long-form answers are annotated. The human-written long-form answers from all the splits of ASQA are written by the same set of annotators using the same annotation template. This creates a formatting bias that models finetuned using thousands of examples can capture. We attribute the low ROUGE-L scores of few-shot systems, in comparison to finetuned systems, to this bias. The second example in Table 8 in Appendix B shows this bias, where both answers are equally answer

| Criteria | FLUENT | DISAMBIG | VS GOLD |
|---------------|--------|----------|---------|
| PaLM best | 0.86 | 0.12 | 0.01 |
| T5 best | 0.80 | -0.17 | 0.04 |
| PaLM few-shot | 0.90 | 0.05 | -0.05 |

Table 6: Human evaluation. Percentage of FLUENT outputs and Best-Worst Scaling scores in terms of DISAMBIGuation preference and similarity VS GOLD answers. Red-colored values are those that are significantly outperformed by at least one other value (paired sample t-test for FLUENT and one-way ANOVA with posthoc Tukey HSD tests for DISAMBIG and VS GOLD, $p < 0.05$)

the question, but are written differently and thus have a significant difference in ROUGE-L. This finding is related to that of Goyal et al. (2022), where they showed that document summaries generated by GPT-3 are always preferred by humans but are ranked the lowest by ROUGE. Future work on evaluation should explore reference-free evaluation, where answers are instead compared to trustworthy sources, similar to Rashkin et al. (2021).

Human Evaluation We conducted human evaluation on 50 randomly selected outputs of PaLM best (prompt tuning using the full dataset), T5 best (open book fine-tuning 5 passages), and PaLM few-shot (100-shot prompt tuning). We asked three expert annotators (1) to mark each output if it has fluency issues and (2) to classify among the outputs which is the best/worst in terms of (2.a) their preference on how the question is disambiguated and (2.b) its similarity to the gold answer. The order of summaries was randomized per participant. For (1), we aggregated the annotations by averaging. For (2) we obtained the best-worst scaling (Louveire et al., 2015) scores, i.e., a rating per system was computed as the percentage of times it was chosen as best minus the percentage of times it was selected as worst. See Appendix C for more details and for an example evaluation template.

Table 6 reports the human evaluation results. As can be seen, there is no significant difference amongst systems in terms of fluency. There is also no significant difference amongst systems in terms of similarity vs gold answers, where T5 is slightly preferred by annotators than both the PaLM-based models. However, T5 outputs are preferred worse in terms of how the question is disambiguated – the best version of PaLM significantly was voted significantly better than T5. These confirm our analysis above. That is, there are other ways to disambiguate the ambiguous question that can actually be preferred more by humans, and T5 copies the summary format of the gold answers to obtain better scores from automatic metrics. Based on our analyses and results from this human elicitation study, we urge the community to rethink the use of gold-standard long-form answers for automatic evaluation.

5 Experiments on Query-focused Multi-document Summarization

5.1 Dataset

The AQuAMuSe dataset⁶ (Kulkarni et al., 2020) is a query-focused multi-document summarization dataset, which was created to simulate how a search engine would consolidate information from multiple documents of high relevance to a given query. The dataset is also a subset of the NQ dataset (Kwiatkowski et al., 2019), but is extended with web documents extracted from Common Crawl and long-form answers from Wikipedia. The dataset consists of 6,599, 714, and 849 training, development, and test examples, where each example is given, on average, 6.46 web documents (2,008 tokens per document). We note that in the closed book setting, the web documents are not used.

5.2 Evaluation

We compare systems using an n-gram overlap-based metric ROUGE-1/2/L (Lin, 2004) and a QA-based metric QAEval_r (Narayan et al., 2022). QAEval_r is similar to the original QAEval (Deutsch et al., 2021), where a question generation model is used to generate questions from the gold summary, and a question answering model attempts to answer these questions using the system-generated summary as the context. In QAEval_r, the way questions are generated is modified such that we only consider questions that are information-seeking (i.e., based on the theme-rheme structure; Vallduví and Vilkkuna, 1998; Kruijff-Korbayová and Steedman, 2003).

5.3 Results

We compared the following systems. Few-shot prompting systems include PaLM 540B using random exemplars, and using NL/QA/AF refinements. For the refinements, we used the 20 exemplars labeled with the *Needs Elaboration* type, which are sampled from the AQuAMuSe training dataset (i.e., the sixth type in Table 2). We also experimented including from ASQA prompt exemplars during dynamic prompt selection to allow the model to more effectively differentiate question types. Few-shot prompt tuning systems include PaLM 540B, with and without AF refinement, and with ASQA prompt exemplars during dynamic prompt selection. Finally, fully finetuned systems include T5-XL closed book and open book, LongT5-XL (Guo et al., 2022) open book that allows longer contexts, and PaLM 540B prompt tuning. For the open book systems, we filled in the context with as many passages as possible within the limits of their input lengths. This results to T5 having 8.3 passages and LongT5 having 30.0 passages on average.

Table 7 reports the scores of these systems. Here, we see similar results as in ASQA; in the few-shot prompting and prompt tuning systems, the addition of AF re-

⁶https://www.tensorflow.org/datasets/community_catalog/huggingface/aquamuse

| System | R1 | R2 | RL | QAE _r |
|------------------------------------|--------------|--------------|--------------|------------------|
| <i>Few-shot Prompting</i> | | | | |
| PaLM 540B (random) | 31.52 | 15.66 | 28.12 | 10.26 |
| + NL refinement | 37.52 | 17.34 | 33.62 | 12.20 |
| + QA refinement | 34.78 | 16.97 | 31.21 | 11.50 |
| + AF refinement | 36.84 | 16.92 | 32.94 | 12.50 |
| + ASQA exemplars | 37.72 | 18.11 | 33.73 | 13.24 |
| <i>Few-shot Prompt Tuning (PT)</i> | | | | |
| PaLM 540B 20-shot | 34.12 | 11.89 | 29.63 | 7.96 |
| + AF refinement | 34.33 | 12.43 | 30.03 | 9.13 |
| + ASQA exemplars | 37.45 | 16.55 | 33.30 | 11.77 |
| <i>Using the full dataset</i> | | | | |
| T5-XL CB | 29.39 | 10.45 | 26.17 | 5.62 |
| T5-XL OB 8.3 Passages | 44.93 | 27.10 | 41.53 | 19.77 |
| LongT5-XL OB 30.0 Passages | 64.11 | 50.60 | 61.43 | 39.22 |
| PaLM 540B PT | 40.53 | 19.66 | 36.56 | 14.62 |

Table 7: Evaluation of several systems on the test set of the AQuAMuSe dataset. The best values for each block are **bold-faced**. CB and OB refer to closed-book and open-book settings.

finements improves the most over the random baseline among the three types of refinements. Moreover, the inclusion of ASQA prompt exemplars during dynamic prompt selection also substantially improves the scores. It is interesting to note that, in contrast to ASQA results, prompt tuning performs worse than few-shot prompting, having worse QAEval_{theme} scores overall. When compared with finetuning systems, our prompting systems significantly outperform the closed book variant of T5, but fall behind the open book systems. We believe that this is due to the fact that AQuAMuSe summaries are extractive – only 2%/13%/24%/31% of the unigrams/bigrams/trigrams/4-grams in the summaries are novel. This allows open-book systems to just copy directly from the source.

6 Conclusions

In this paper, we investigated the ability of large language models to answer questions in a long-form manner through two question answering datasets: ASQA and AQuAMuSe. We introduced query refinement prompting that improves over standard few-shot prompting and prompt tuning methods by encouraging the model to explicitly express the multifacetedness in questions. With the use of query refinement prompts on both few-shot closed book prompting and prompt tuning settings, we are able to outperform systems trained using the full training data in both ASQA and AQuAMuSe. While we also achieve comparable results with open book systems in ASQA, we acknowledge that these systems still perform better when they have access to more retrieved passages. For future work, we plan to explore ways to few-shot prompt large language models in the open book setting and ways to augment a retrieval component into large language models.

7 Limitations

There are two major limitations in this work. Firstly, while we showed that query refinement prompts improve the ability of LLMs to generate long-form answers in a closed book and few-shot settings, open-book systems still perform better even when using a lot less parameters. Doing open-book long-form question answering in LLMs is currently not trivial due to their input token length limit and the need to use longer prompts when context passages are given.

The other major limitation is that human annotators only used the gold-standard answers to check for correctness (i.e., VS GOLD in Table 6). As explained in Section 4.5, there can be many ways to disambiguate questions, and therefore the systems can obtain a long-form answer that is not different from the gold answers but still should be considered correct. We tried asking annotators to use the Internet to check for correctness, however they found it difficult to do so even for a single example.

Acknowledgements

The authors would like to thank Tal Schuster and Mirella Lapata for their detailed feedback, and Ji Ma and Priyanka Agrawal for their help in prompt tuning experiments. Finally, we are thankful to the PaLM team in Google Research for providing easy-to-use tools to experiment with large language models.

References

- Ankur Bapna and Orhan Firat. 2019. [Simple, scalable adaptation for neural machine translation](#). 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 1538–1548, Hong Kong, China. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentaoh Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. [QuAC: Question answering in context](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2174–2184, Brussels, Belgium. Association for Computational Linguistics.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *arXiv preprint arXiv:2204.02311*.
- Jordan Clive, Kris Cao, and Marek Rei. 2021. [Control prefixes for text generation](#). *arXiv preprint arXiv:2110.08329*.
- Hoa Trang Dang. 2005. Overview of duc 2005. In *Proceedings of the document understanding conference*, volume 2005, pages 1–12.
- Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. [Towards question-answering as an automatic metric for evaluating the content quality of a summary](#). *Transactions of the Association for Computational Linguistics*, 9:774–789.
- Katherine Elkins and Jon Chun. 2020. [Can gpt-3 pass a writer’s turing test?](#) *Journal of Cultural Analytics*, 5(2):17212.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. [ELI5: Long form question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3558–3567, Florence, Italy. Association for Computational Linguistics.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. [News summarization and evaluation in the era of gpt-3](#). *arXiv preprint arXiv:2209.12356*.
- Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. [LongT5: Efficient text-to-text transformer for long sequences](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 724–736, Seattle, United States. Association for Computational Linguistics.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. [Towards a unified view of parameter-efficient transfer learning](#). In *International Conference on Learning Representations*.
- Haozhe Ji, Pei Ke, Shaohan Huang, Furu Wei, and Minlie Huang. 2020. [Generating commonsense explanation by extracting bridge concepts from reasoning paths](#). In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 248–257, Suzhou, China. Association for Computational Linguistics.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of*

- the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*.
- Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. [Hurdles to progress in long-form question answering](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4940–4957, Online. Association for Computational Linguistics.
- Ivana Kruijff-Korbayová and Mark Steedman. 2003. Discourse and information structure. *Journal of logic, language and information*, 12(3):249–259.
- Sayali Kulkarni, Sheide Chammas, Wan Zhu, Fei Sha, and Eugene Ie. 2020. Aquamuse: Automatically generating datasets for query-based multi-document summarization. *arXiv preprint arXiv:2010.12694*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). 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 4582–4597, Online. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Xiaochen Liu, Yu Bai, Jiawei Li, Yinan Hu, and Yang Gao. 2022. Psp: Pre-trained soft prompts for few-shot abstractive summarization. *arXiv preprint arXiv:2204.04413*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Jordan J Louviere, Terry N Flynn, and Anthony Alfred John Marley. 2015. *Best-worst scaling: Theory, methods and applications*. Cambridge University Press.
- Sewon Min, Kenton Lee, Ming-Wei Chang, Kristina Toutanova, and Hannaneh Hajishirzi. 2021. [Joint passage ranking for diverse multi-answer retrieval](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6997–7008, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. [AmbigQA: Answering ambiguous open-domain questions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5783–5797, Online. Association for Computational Linguistics.
- Shashi Narayan, Joshua Maynez, Reinald Kim Amplayo, Kuzman Ganchev, Annie Louis, Fantine Huot, Dipanjan Das, and Mirella Lapata. 2022. Conditional generation with a question-answering blueprint. *arXiv preprint arXiv:2207.00397*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. 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.
- Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. 2021. Measuring attribution in natural language generation models. *arXiv preprint arXiv:2112.12870*.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. [CoQA: A conversational question answering challenge](#). *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. [How much knowledge can you pack into the parameters of a language model?](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426, Online. Association for Computational Linguistics.

- Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp*, 109:109.
- Pedro Rodriguez, Shi Feng, Mohit Iyyer, He He, and Jordan Boyd-Graber. 2019. [Quizbowl: The case for incremental question answering](#).
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. [Learning to retrieve prompts for in-context learning](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M. Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. 2021. [Multitask prompted training enables zero-shot task generalization](#). *CoRR*, abs/2110.08207.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. [BLEURT: Learning robust metrics for text generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- Charlie Snell, Dan Klein, and Ruiqi Zhong. 2022. [Learning by distilling context](#). *arXiv preprint arXiv:2209.15189*.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. [Asqa: Factoid questions meet long-form answers](#). *arXiv preprint arXiv:2204.06092*.
- Enric Vallduví and Maria Vilkkuna. 1998. On rheme and kontrast. In *The limits of syntax*, pages 79–108. Brill.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. [Finetuned language models are zero-shot learners](#). In *International Conference on Learning Representations*.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. [Finetuned language models are zero-shot learners](#). *CoRR*, abs/2109.01652.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022b. [Chain of thought prompting elicits reasoning in large language models](#). *arXiv preprint arXiv:2201.11903*.
- Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. 2021. [QMSum: A new benchmark for query-based multi-domain meeting summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5905–5921, Online. Association for Computational Linguistics.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed Chi. 2022. [Least-to-most prompting enables complex reasoning in large language models](#). *arXiv preprint arXiv:2205.10625*.

A Prompt Templates

This section presents the prompt templates we used in our experiments. Figures 1 and 2 show prompt templates with AF refinements for ASQA and AQuA-MuSe, respectively. When using no refinements, the Disambiguations part of the template is removed. When using a different query refinement step, the same part is replaced with a natural language explanation or a list of question answers. Finally, Figure 3 shows the template used when the oracle disambiguated questions are given.

I will provide ambiguous questions that have multiple answers about different aspects of the question, and answer them in detail with at least two sentences.

Question: Who sang it's a long way to the top?
Disambiguations:
- band: AC/DC
- lead vocal: Bon Scott
Answer: "It's a Long Way to the Top (If You Wanna Rock 'n' Roll)" is a song by Australian hard rock band AC/DC. This was a signature song for lead singer Bon Scott. Brian Johnson, who replaced Scott as AC/DC's lead vocalist after Scott's death in 1980, does not perform it, out of respect for his predecessor.

Question:

Figure 1: Example ASQA Prompt with AF refinements and one exemplar.

B Annotator Bias Examples

Table 8 shows two sets of examples from the ASQA dataset where PaLM answered the question well but received a lower ROUGE-L score than T5. The first example shows that there can be multiple correct ways to disambiguate and answer the question. In the example, the gold and the T5 answers disambiguated the question based on the habitat of the animal (fresh water or salt-water), while the PaLM answer disambiguated it based on both the habitat and the characteristics (having gills, having lungs, or having both).

The second example shows that T5 generates answers that follow a dataset-specific format. In the example, while both systems correctly answered *Tony Hawk* and *Danny Way*, the T5 long answer looks more similar to the gold-standard answer, whereas the PaLM long answer looks more concise and arguably more readable.

C Human Evaluation Instructions and Template

In this section, we share the instructions and template used for human evaluation. The annotators did the annotation as part of their jobs, and consent was given

I will provide questions that need to be elaborated to be answered fully, and will answer them in detail with at least two sentences.

Question: where did the term shooting brake come from

Details:

- how the term originated: as an early 19th century British term

- what it was for: a vehicle used to carry shooting parties with their equipment and game

- etymology of the term brake: uncertain; initially a chassis used to break in horses, used to describe a motorized vehicle

- its possible origins: in the Dutch word 'brik' which means 'cart' or 'carriage' "

Answer: "Shooting-brake" originated as an early 19th century British term for a vehicle used to carry shooting parties with their equipment and game. The etymology of the term brake is uncertain; initially a chassis used to break in horses, and subsequently used to describe a motorized vehicle. It is also possible, that the word 'brake' has its origins in the Dutch word 'brik' which means 'cart' or 'carriage'.

Question:

Figure 2: Example AQuAMuSe Prompt with AF refinements and one exemplar.

I will provide ambiguous questions that have multiple answers about different aspects of the question, and answer them in detail with at least two sentences.

Question: Who sang it's a long way to the top?

Disambiguated Questions:

Q: Which band sang it's a long way to the top?

Q: Who was the lead vocal of it's a long way to the top?

Disambiguated Answers:

- band: AC/DC

- lead vocal: Bon Scott

Answer: "It's a Long Way to the Top (If You Wanna Rock 'n' Roll)" is a song by Australian hard rock band AC/DC. This was a signature song for lead singer Bon Scott. Brian Johnson, who replaced Scott as AC/DC's lead vocalist after Scott's death in 1980, does not perform it, out of respect for his predecessor.

Question:

Figure 3: Example ASQA Prompt with AF refinements, oracle disambiguated questions, and one exemplar.

Question: What do you call animals live in water?

Gold: **Marine life, or sea life or ocean life**, is the plants, animals and other organisms that live in the salt water of the sea or ocean, or the brackish water of coastal estuaries. An **aquatic animal** is an animal, either vertebrate or invertebrate, which lives in the water for most or all of its lifetime. The term aquatic can be applied to animals that live in either fresh water or salt water. However, the adjective marine is most commonly used for animals that live in saltwater, i.e. in oceans, seas, etc.

T5 best: **aquatic animals** are animals that live in water for most or all of their lifetime. aquatic animals can be either vertebrate or invertebrate. fresh water animals are **aquatic animals**. **marine animals** are animals that live in saltwater. **ROUGE-L:** 44.96

PaLM best: Animals that live in water are called **aquatic animals**. Aquatic animals that have gills are called **fish**. Aquatic animals that have lungs are called **marine mammals**. Aquatic animals that have both gills and lungs are called **amphibians** **ROUGE-L:** 20.16

Question: Who landed the first 900 on a skateboard?

Gold: The 900 is a 2 1/2-revolution (900 degrees) aerial spin performed on a skateboard ramp. A number of opinions exist about pre-1999 900s. The most prominent of these is the argument that **Danny Way** landed the 900 in 1989, and it appears in an early Santa Cruz film. **Tony Hawk**, one of the most successful vertical pro skateboarders in the world, is generally said to have landed the first 900 on a skateboard. He landed the 900 at X Games V in 1999 after ten failed attempts.

T5 best: the 900 is a 2 1/2-revolution (900 degrees) aerial spin performed on a skateboard ramp. a number of opinions exist about pre-1999 900s. the most prominent of these is the argument that **danny way** landed the 900 in 1989, and it appears in an early santa cruz film. in 1999 **tony hawk**, one of the most successful vertical pro skateboarders in the world, landed "the 900" at the 1999 x-games after ten failed attempts. in october 2011, schaar became the 8th person in history to land a 900 on a skateboard. in september 2016, **schaar** became the 1st person in history to land a stalefish 900 on a skateboard. **ROUGE-L:** 80.00

PaLM best: The 900 is a skateboarding trick in which the skateboarder does two and a half vertical revolutions. American professional skateboarder **Tony Hawk** landed the first 900 in a competition at the X Games in 1999. Australian professional skateboarder **Danny Way** landed the first 900 on video in 1990. **ROUGE-L:** 45.59

Table 8: ASQA outputs from the T5 and PaLM best configurations. Disambiguated answers are **bold-faced**.

by all annotators that the annotations will be used as part of evaluation. They are given a spreadsheet to conduct the annotation process, and is given the following instruction:

You are given a question and three system-generated answers. Please read them carefully and do the following tasks:

1. **Fluency:** For each system-generated answer, check the box if you find fluency errors in the answer. That is, answers with grammatical, spelling, and coherence issues should have their boxes checked.
2. **Disambiguation:** Among the system-generated answers, select the best/worst answer based on following criterion: The system-generated answer contains different answers to the question which coincide with my information need. You are not required to check for correctness at this stage.
3. **Similarity with Gold:** Given the following gold-standard human-written answer to the question, select the best/worst answer among the system-generated answer based on the following criterion: Based on your overall impression, the system-generated answer contains and only contains all the information shown in the human-written answer.

Figure 4 resembles how the spreadsheet looks like. Annotators are asked to do the tasks from left to right, starting with the fluency error checkboxes, then the disambiguation best/worst summary selection, then finally the similarity with gold best/worst summary selection. By doing this, their judgments on the first two tasks are not influenced by the gold answer.

| Question | Answer1 | FE? | Answer2 | FE? | Answer3 | FE? |
|---|--|--------------------------|---|--------------------------|---|--------------------------|
| What family does paris belong to in romeo and juliet? | In William Shakespeare's Romeo and Juliet, the Capulets are one of the two feuding families in Verona, the other being the Montagues. The Capulets are represented by Capulet, his wife Lady Capulet, their son Tybalt, their daughter Juliet, and Juliet's Nurse. | <input type="checkbox"/> | Count Paris or County Paris is a fictional character in William Shakespeare's Romeo and Juliet. He is a suitor of Juliet. He is handsome, wealthy, and a kinsman to Prince Escalus. In the play, Paris is portrayed as the patriarch of the Capulet family, the Capulet family being the name of an actual political faction from the 13th century. | <input type="checkbox"/> | In William Shakespeare's Romeo and Juliet, Paris is a young nobleman, kinsman to Prince Escalus, and suitor to Juliet. Paris is a Capulet, a member of the House of Capulet, one of the two feuding families in Verona. Paris is a Montague, a member of the House of Montague, the other feuding family in Verona. | <input type="checkbox"/> |

| Disambiguation | | Gold | Similarity with Gold | |
|----------------|-------|--|----------------------|-------|
| Best | Worst | | Best | Worst |
| ○ 1 | ○ 1 | Count Paris is a kinsman of Prince Escalus and seeks to marry Juliet. He belongs to the House of Escalus which also includes Prince Escalus, the Prince of Verona, who is the desperate resolver of the feuding families and Mercutio, who is the cousin of Prince Escalus and Count Paris, and is a close friend of Romeo and his cousin Benvolio. He supports and fights on the Montague side of the feud, and just like a Montague, hates the Capulet family. | ○ 1 | ○ 1 |
| ○ 2 | ○ 2 | | ○ 2 | ○ 2 |
| ○ 3 | ○ 3 | | ○ 3 | ○ 3 |

Figure 4: Human evaluation template with one example. The table below is connected to the table above from left to right.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Section 7
- A2. Did you discuss any potential risks of your work?
There is no potential risk as we do not plan to release models.
- A3. Do the abstract and introduction summarize the paper’s main claims?
Section 1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Sections 4.1 and 5.1

- B1. Did you cite the creators of artifacts you used?
Sections 4.1 and 5.1
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
The license and terms can be found in the link provided in Sections 4.1 and 5.1.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Sections 4.1 and 5.1
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Not applicable. Left blank.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Not applicable. Left blank.
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Sections 4.1 and 5.1

C Did you run computational experiments?

Sections 4 and 5

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Sections 4 and 5

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
Sections 4 and 5
- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?
Sections 4 and 5
- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?
Sections 4 and 5
- D** **Did you use human annotators (e.g., crowdworkers) or research with human participants?**
Sections 4
- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?
Appendix C
- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?
Appendix C
- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?
Appendix C
- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?
Not applicable. Left blank.
- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?
Since we only used three expert annotators, we do not want to share these information as they might be identified using such information.