

Crafting In-context Examples according to LMs’ Parametric Knowledge

Anonymous ACL submission

Abstract

In-context learning has been applied to knowledge-rich tasks such as question answering. In such scenarios, in-context examples are used to trigger a behaviour in the language model: namely, it should surface information stored in its parametric knowledge. We study the construction of in-context example sets, with a focus on the parametric knowledge of the model regarding in-context examples. We identify ‘known’ examples, where models can correctly answer from its parametric knowledge, and ‘unknown’ ones. Our experiments show that prompting with ‘unknown’ examples decreases the performance, potentially as it encourages hallucination rather than searching its parametric knowledge. Constructing an in-context example set that presents both known and unknown information performs the best across diverse settings. We perform analysis on three multi-answer question answering datasets, which allows us to further study answer set ordering strategies based on the LM’s knowledge about each answer. Together, our study sheds lights on how to best construct in-context example sets for knowledge-rich tasks.

1 Introduction

Large language models (LLMs) can perform competitively on knowledge-rich tasks such as question answering via in-context demonstrations (Brown et al., 2020). In such scenarios, in-context examples are used not only to teach the LLM the mapping from inputs to outputs, but also to invoke the LLM’s parametric knowledge (Liu et al., 2021; Agrawal et al., 2022). Given such role of in-context examples, we examine how the LLM’s parametric knowledge of in-context examples impact the effectiveness of in-context examples.

Let’s imagine a very challenging in-context example set, where LLMs cannot answer any of in-context examples from its parametric knowledge.

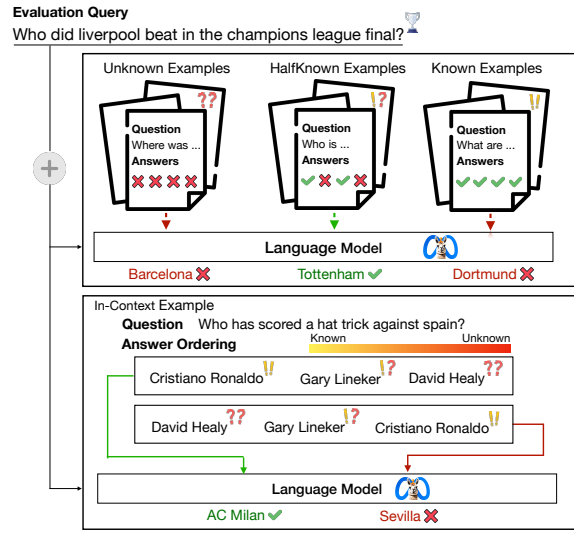


Figure 1: We study how an LM’s knowledge of in-context examples impacts their effectiveness. On the top box, we construct three sets of in-context examples, Unknown, HalfKnown, and Known, differing in its difficulty (Section 3). On the bottom box, we construct two in-context examples, which contain the same question and answer set but answers are sorted differently: one in increasing amount of parametric knowledge and one in reverse (Section 4, 5).

For example, in-context examples can query knowledge about recent events that happened after pre-training. These in-context examples will teach the model to generate plausible-looking responses, but may encourage hallucination as a result. On the other hand, if we only provide in-context examples where LLM can easily answer, would LLM learn to make an educated guess on more challenging evaluation examples?

We pose a suite of research questions connecting parametric knowledge of an LM on in-context examples and its impact on model predictions. Figure 1 provides our study overview. We mainly evaluate on multi-answer QA datasets (Min et al., 2020; Malaviya et al., 2023; Amouyal et al., 2022), a challenging knowledge-rich task, and a math QA

dataset (Cobbe et al., 2021), which requires reasoning from LLM. Multi-answer QA datasets further allows a controlled study where we fix the question and vary a choice of answer from a set of valid answers, or how we order answers based on model’s parametric knowledge of individual answer.

We first compare providing ‘known’ or ‘unknown’ in-context examples (Section 3). We operationalize ‘known’ in-context examples as those LM can correctly predict with in-context learning. We do not observe a clear winner between two choices, with results varying depending on the dataset. Throughout all datasets, however, providing in-context examples that have a mixture of known and unknown information leads to superior performance compared to solely known or unknown in-context examples.

Our next analysis focuses on the ordering of multi-answer set while fixing in-context example set (Section 4, 5). Compared to randomly ordering valid answers, semantically meaningful ordering brings substantial changes in model predictions. Even alphabetical ordering of answer set changes predicted answers substantially, prompting model to generate 1.5 more answer on average than when shown randomly sorted answer set. We further find that ordering the answer set of in-context examples in descending order of model knowledge often leads to performance gains. Together, our work suggests best practices for crafting in-context examples, with relation to its parametric knowledge, for knowledge-intensive tasks.

2 Experimental Settings

We first describe our evaluation setting which centers around multi-answer QA datasets.

2.1 Dataset

We evaluate on three multi-answer QA datasets: (1) AmbigQA (Min et al., 2020) contains a subset of questions from the Natural Questions (Kwiatkowski et al., 2019) dataset, namely those marked as ambiguous in the sense that depending on the interpretation they can have multiple correct answers. (2) QAMPARI (Amouyal et al., 2022) consists of questions whose set of correct answers necessarily span multiple paragraphs in the document from which they were retrieved. The dataset was originally developed to evaluating retrieval methods, and we repurpose it to create a challenging closed-book QA setting.

(3) QUEST (Malaviya et al., 2023) dataset is constructed by formulating queries that define implicit set operations over Wikipedia entities. We report the dataset statistics in Appendix A.

2.2 Evaluation Metrics

Given a question q , the model will predict a set of answers $\hat{a} = \{a_1, a_2, \dots, a_m\}$, where each $a_i = (w_{i_1}, w_{i_2}, \dots, w_{i_{|a_i|}})$ is a sequence of tokens for a single answer. We denote $a^* = \{a_1^*, a_2^*, \dots, a_n^*\}$ as the ground truth answers to the same question.

We use standard token match metrics for evaluating answer accuracy, Exact Match (EM) and F1-score (Joshi et al., 2017). EM assigns a score of 1 if the predicted answer equals to the ground truth answer, while F1-score is calculated over the tokens in the answer. We use metrics for multi-answers introduced in prior work (Min et al., 2020), which we describe below for completeness.

Answer-level Exact Match ($F1_{EM}$) As predicting the exact ground truth answer set correctly is very challenging, we report the F1-score of answer-level exact match, denoted as $F1_{EM}$. For an answer a and reference answers set S , we define a correctness score $c(f, a, S) = f(a, S)$ with respect to function f . We use $f(a, S) = \mathbb{1}(a \in S)$ here. Then, we calculate the F1-score over set-level precision and recall according to c .

$$P = \frac{\sum_{i=1}^m c(f, a_i, a^*)}{m}, R = \frac{\sum_{j=1}^n c(f, a_j^*, \hat{a})}{n}$$

$$F1_{EM} = \frac{2 \times P \times R}{P + R}$$

Answer-level F1 ($F1_{F1}$) The generated answer may be semantically equivalent to one of the ground truth answers, without being lexically equivalent (e.g., "Friends" and "The TV show Friends"). To account for such semantic equivalences, we use $F1$ score between the tokens of two answer strings instead of the exact match as a correctness score, $f(a, S) = \max_{a' \in S} (F1(a, a'))$. Then, we compute F1-score over set-level precision and recall as above.

Statistical Testing As our evaluation datasets are relatively small, we conduct paired bootstrap tests throughout most of our experiments, highlighting results that outperform baseline with p value of ≤ 0.05 .

	AmbigQA _{dev}		QAMPARI _{dev}		QUEST _{test}	
	Llama2	GPT-3.5	Llama2	GPT-3.5	Llama2	GPT-3.5
Random	18.0 / 28.9	20.0 / 31.6	10.3 / 20.8	15.0 / 28.5	3.4 / 11.0	6.0 / 16.6
Unknown	17.2* / 28.2*	20.3* / 33.1*	10.9* / 22.0*	14.8 / 27.9*	3.7* / 11.9*	5.7* / 15.8*
HalfKnown	18.5* / 29.5*	21.6* / 33.2*	11.3* / 22.6	15.5* / 28.2*	4.0* / 11.9*	6.3* / 17.4*
Known	18.3* / 29.0*	21.3* / 33.1*	9.8 / 19.7	15.3 / 29.2*	3.9* / 12.0*	5.4* / 15.8

Table 1: Results comparing known example and unknown example. We present $F1_{EM}$ and then $F1_{F1}$ in each cell. Using half-known example outperforms other settings. We put * on scores that are significantly different from that of Random in-context examples set, and bold the highest performing set for each metric.

2.3 Base Models

Language Model We evaluate on Llama2 (Touvron et al., 2023) (13B) language model mainly and additionally OPT (Zhang et al., 2022) (13B) and GPT-3.5-turbo models to evaluate generalization.

In-context Example Retriever Prior work (Rubin et al., 2021) has established that using semantically similar in-context examples improves the performance of in-context learning significantly. Throughout our study, we often retrieve top 5 most similar in-context examples from the entire training set for each dataset to form the prompt. We place in-context examples in decreasing order of similarity, such that the most similar example will be presented closest to the evaluation question. We measure example similarities by encoding each question with a SimCSE model (Gao et al., 2021) and computing their dot product.

3 Known Examples vs. Unknown Examples

Prior work has studied a few characteristics of successful in-context example set, such as label distribution in the in-context example set (Min et al., 2022). We evaluate in-context examples with respect to model’s parametric knowledge, whether a “known” or “unknown” in-context example is better. We operationalize “known” ones as the ones where LLMs can get the answers correctly from its own parametric knowledge, and “unknown” ones as those that cannot be answerable from its parametric knowledge.

3.1 In-context Example Set Study

We create four sets of in-context examples, differing in its difficulty for a given LM.

- **UNKNOWN**: examples for which the LM possesses no knowledge of the answers. Operationally, these are examples when LM is

prompted with five most similar examples, LM will predict zero answer correctly (i.e. zero $F1_{EM}$ score).

- **RANDOM**: randomly sampled examples. Since the LM possesses no knowledge to majority of the examples, these exhibit 0.18 $F1_{EM}$ score on average.
- **HALFKNOWN**: examples for which the LM possesses roughly half knowledge of the answers (i.e. 0.5 $F1_{EM}$ score).
- **KNOWN**: examples for which the LM possesses full knowledge of the answers (i.e. 1.0 $F1_{EM}$ score).

As prior work (Rubin et al., 2021) has established that the similarity of in-context example to the query correlates strongly with the model’s performance, we control for this confounding factor. We compute the average similarity for each in-context example candidate to other in-context example candidates in the candidate set (training set). Then, we choose a fixed number of in-context examples whose average similarity value is close to the median value.¹ From this candidate set, we sample five examples for each condition and use them as fixed in-context examples across all questions in the evaluation dataset. To further reduce randomness, we sample multiple sets of five example set for each condition and report the average performance (by default, four sets are sampled and two sets are sampled for HALFKNOWN and KNOWN set in QUEST because of lack of examples with sufficient model knowledge).

We present the performance of each in-context example set for three datasets with Llama2 and GPT-3.5 in Table 1. We observe the HALFKNOWN

¹We choose 999 examples for AmbigQA and QAMPARI, and 499 for QUEST (as QUEST only has 1251 training examples), half from below median, half from above median. For QUEST, we could not find enough examples with where model score full $F1_{EM}$ score, so we selected highest scoring examples. The mid-range is (0.245, 0.264), (0.294, 0.296), (0.326, 0.373) for AmbigQA, QAMPARI, and QUEST.

Unknown	Random	HalfKnown	Known
33.1	34.8	36.4	32.0

Table 2: The accuracy on GSM8K dataset. Accuracy is expressed as the percentage of correct answers over the entire test dataset, which consists of 1319 queries.

in-context example set achieves strong performance consistently on both LMs. Since HALF-KNOWN with in-context examples that contain both answers that the model knows and doesn’t know, we hypothesize this may successfully prompt LMs to leverage parametric knowledge and to make educated guesses.

3.2 Analysis

In this section, we provide two additional studies with Llama2 model.

Extension to Math QA Dataset We explore constructing in-context example sets with varying “knowness” for single-answer QA task. We chose GSM8K (Cobbe et al., 2021) dataset, a commonly used dataset for investigating the reasoning capabilities of LLMs. GSM8K consists of 8,500 natural language questions requiring arithmetic reasoning for obtaining an answer. To evaluate parametric knowledge available to solve each training example with the LM, we prompt each example with the 8-shot example set taken from Wei et al. (2022b) and classified as correct, wrong, or invalid, where invalid indicates that the model did not produce an answer. We construct four in-context example sets:

- UNKNOWN set includes randomly selected six examples that model answered incorrectly.
- RANDOM set includes randomly selected six examples from entire training dataset. The LM correctly answer questions in training set for 20% of questions.
- HALFKNOWN set includes three correct and three wrong examples.
- KNOWN set includes randomly selected six examples that model answered correctly.

We select six examples four times and report the averaged accuracy in Table 2. HALFKNOWN set achieves the highest accuracy, repeating the trend from multi-answer QA datasets.

Single Answer Study In this study, we further control for variability in the question used in in-context examplars. We fix the in-context example

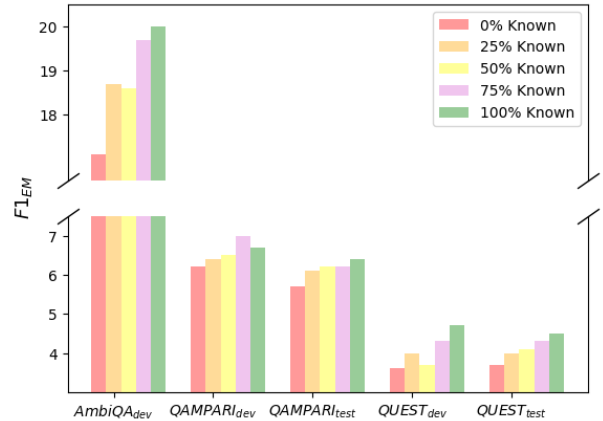


Figure 2: Results of single answer study on Llama2 model. Only an answer at the x -th quantile of perplexities in decreasing order is presented in each in-context example. As the model gets exposed to more known answers, the performance tend to increase.

set and manipulate the multi-answer set, such that we provide only one answer from multi answer set for each in-context example. For example, if a question in in-context example is “who was the president of U.S.?”, we can either provide a famous president or a lesser-known president as an answer. Both are “correct” answers, but which answer would lead to better model performance?

For each question in our evaluation set, we retrieve top five most similar examples in training set as in-context examples. We will measure perplexity of each answer to approximate how well LM ‘knows’ the answer. For each example, a pair of question q and gold answer set $\{a_1^*, a_2^*, \dots, a_n^*\}$, we form a prefix p by prepending top five most similar examples to the query q .² Then, we compute the length normalized perplexity of each answer a_i^* and prefix p as follows:

$$PP(a_i^*|p) = \prod_{j=1}^{|a_i^*|} P(w_{ij}|p, w_{i1}, \dots, w_{i(j-1)})^{-\frac{1}{|a_i^*|}} \quad 282$$

We will order the gold answer set in descending order of perplexity, and select an answer at the x -th quantile. This way, an answer at the 100% quantile represents the most ‘known’ answer, as its perplexity is the lowest among the gold answers.

Figure 2 presents the $F1_{EM}$ score among various x -th quantile. We observe a clear trend across all three datasets, that using a ‘known’ answer leads LM to generate more accurate answer. These in-context examples are incomplete, only presenting

²We present an example prefix in Appendix G.

one answer while there are multiple valid answers. This leads to low performance overall, as LM will only generate a single answer (low recall). Yet, this experiment affirms that crafting in-context example by considering model’s parametric knowledge can impact the final performances.

4 Ordering Answers Based on LM’s Knowledge

Prior work suggests that the ordering of in-context examples significantly impacts the performance, with more relevant examples being most beneficial when placed last (Zhao et al., 2021). Yet, no prior work has studied the ordering of answers inside each in-context example. We investigate this here. Following our previous study, our focus is on **parametric** knowledge of LMs being prompted. Specifically, we question whether placing answers based on how well the model knows about answers improves the performance.

We present strategies to order the answer set of each example, a pair of question q and its gold answer set $a^* = \{a_1^*, a_2^*, \dots, a_n^*\}$, which will be used as an in-context example.³ We present two baselines and two methods (PERPLEXITY, GREEDY) for ordering the gold answer set of each in-context example based on model’s parametric knowledge.

Baselines The RANDOM baseline randomly orders answers, and ALPHABET orders answers alphabetically. While alphabetical ordering is not relevant to model’s parametric knowledge of the answer, prior work (Madaan et al., 2022) has shown that consistent ordering of labels can improve the performance of fine-tuned LLM’s predictions.

Knowledge-Aware Ordering We decide ordering based on the **perplexity** of individual answer given the prefix, or by performing greedy **constrained decoding** given the prefix. We use the same prefix as in Section 3.2, a concatenation of five in-context examples. Each ordering strategy will yield two orderings of answers, which either sorts the answers in the descending order of model’s parametric knowledge or ascending order (denoted as REVERSE).

- **PERPLEXITY:** We compute the length normalized perplexity of each answer a_i^* and prefix p as

³As reordering process is computationally expensive, proportional to the number of answers, we only consider examples that have less than 20 answers. This results in exclusion of 1 example in AmbigQA, 8094 examples in QAMPARI, and none in QUEST.

Input: LM \mathcal{M} , Prefix p , Gold answer set $a^* = \{a_1^*, \dots, a_n^*\}$, where each gold answer is a token sequence (i.e., $a_i^* = (w_{i_1}, \dots, w_{i_{|a_i^*|}})$)

Output: Ordered answer indices of the gold answer set

```

1:  $I_1 \leftarrow \{w_{1_1}, \dots, w_{n_1}\}$ 
2:  $u \leftarrow 1$ 
3: while  $I_1 \neq \emptyset$  do
4:    $t \leftarrow 0$ 
5:   repeat
6:      $t \leftarrow t + 1$ 
7:      $o_t \leftarrow \operatorname{argmax}_{w \in I_t} P_{\mathcal{M}}(w|p)$ 
8:      $p \leftarrow [p; o_t]$ 
9:      $I_{t+1} \leftarrow \{w_{i_{t+1}} | w_{i_t} == o_t\}$ 
10:  until  $\exists a_{k_u}^* == (o_1, \dots, o_t)$  {this assigns  $k_u$  the index of completed answer}
11:   $I_1 \leftarrow I_1 \setminus \{w_{k_u}\}$ 
12:   $u \leftarrow u + 1$ 
13: return  $\{k_1, \dots, k_n\}$ 

```

Figure 3: Algorithm for constrained decoding for GREEDY ordering.

used in Section 3.2. Then, we sort the answers in ascending order of these perplexities, resulting in ‘known’ answers placed earlier.

- **GREEDY:** We arrange the gold answers by performing a beam search decoding in a greedy manner, constrained to permissible tokens. There will be two loops, outer loop for selecting the first token of the generated answer, and inner loop for completing the chosen first token.

Figure 3 presents the pseudocode, which we explain below. Let’s denote a_i^* as a sequence of tokens $(w_{i_1}, w_{i_2}, \dots, w_{i_{|a_i^*|}})$ for the i -th answer. At each decoding step t , a set of permissible tokens I_t is constructed. Initially, $I_1 = \{w_{1_1}, w_{2_1}, \dots, w_{n_1}\}$, a set of the first token for each potential answer. We choose a token from this set that has the highest likelihood given the prompt, i.e., $o_1 = \operatorname{argmax}_{w \in I_1} P(w|p)$. Then, we update the prefix $p \leftarrow [p; o_1]$. This initiates the inner loop, setting $I_2 = \{w_{i_2} | w_{i_1} == o_1\}$ as a set of second token of answers who starts with the selected first token. This continues until one of the answers a_{k_1} is fully generated. Afterwards, we come back to the outer loop, and the initial set of permissible tokens is set to be $I_1 = \{w_{1_1}, w_{2_1}, \dots, w_{n_1}\} \setminus \{w_{t_1}\}$ excluding $a_{k_1}^*$ which has been already generated. This process continues until all answers has been generated, with a time complexity of $O(n|a_i^*|)$.

5 Results for Answer Ordering Strategies

Having introduced strategies for ordering answers for in-context examples, we study how this im-

		S				
		GREEDY	REVERSE GREEDY	PERPLEXITY	REVERSE PERPLEXITY	ALPHABET
\mathcal{D}_e	AmbigQA _{dev}	71.7 / 66.0	39.2 / 37.2	69.5 / 65.8	38.1 / 34.2	87.4 / 55.6
	QAMPARI _{dev}	69.6 / 60.0	42.2 / 41.0	58.1 / 54.1	46.3 / 45.9	95.0 / 58.9
	QAMPARI _{test}	70.0 / 65.7	43.0 / 41.7	58.8 / 55.8	45.0 / 44.2	94.9 / 58.1
	QUEST _{dev}	78.4 / 63.9	47.2 / 45.8	57.1 / 51.5	49.3 / 48.5	95.7 / 52.1
	QUEST _{test}	81.0 / 63.3	45.7 / 45.3	57.6 / 52.5	48.8 / 47.5	95.6 / 50.8
	Average	74.1	43.5	60.2	45.5	93.7

Table 3: Percentage of generated answer ordering matching in-context examples answer ordering, where we use Llama2 for \mathcal{M} . In each cell, we present the percentage from using corresponding answer ordering strategy first ($\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M})$) and the percentage for randomly ordering answers for control ($\phi(S, \mathcal{D}_t^{S_{\text{random}}}, \mathcal{D}_e, \mathcal{M})$).

pacts the generation of answers with Llama2 and OPT. We first evaluate whether the generated answers mimic the ordering of answers in in-context examples. Then, we evaluate whether the ordering impacts the size and the accuracy of predicted answer set. We also report whether two model’s parametric knowledges are in sync, meaning, if one model knows about one fact, does the other model likely to know the same fact? We overall observe such patterns, particularly for QUEST dataset.

5.1 Does the predicted answer set follow the ordering of in-context answer set?

Throughout in-context learning, the model is expected to learn the pattern shown in the demonstrations. We assess the generated answers to observe whether the model has followed the particular ordering shown in the in-context examples.

Metric We introduce a metric $\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M})$. This measures how much LM \mathcal{M} follows the answer ordering strategy S on evaluation dataset \mathcal{D}_e when using in-context examples from training dataset \mathcal{D}_t whose answered are ordered according to S^t .⁴ When S matches S^t , this metric will measure how much predicted outputs mimic the answer ordering strategy of in-context examples.

Let’s denote $\hat{a}_i = \{a_{i_1}, a_{i_2}, \dots, a_{i_m}\}$ be the list of predicted m answers for the i -th example of an evaluation dataset \mathcal{D}_e , following its generation order from model \mathcal{M} . We reorder the predicted answers from \hat{a}_i with respect to S and denote $f(a_{ij})$ to be the index of a_{ij} in the newly ordered set.

For each consecutive answer pair in \hat{a}_i , we evaluate whether their order is preserved after reordering. Then we count the number of consecutive answer pairs that have preserved the ordering, which is $P_i = \sum_{j=1}^{m-1} \mathbb{1}(f(a_{ij}) < f(a_{i(j+1)}))$. Similarly,

⁴We assume retrieving five most similar in-context examples for each evaluation example throughout this study.

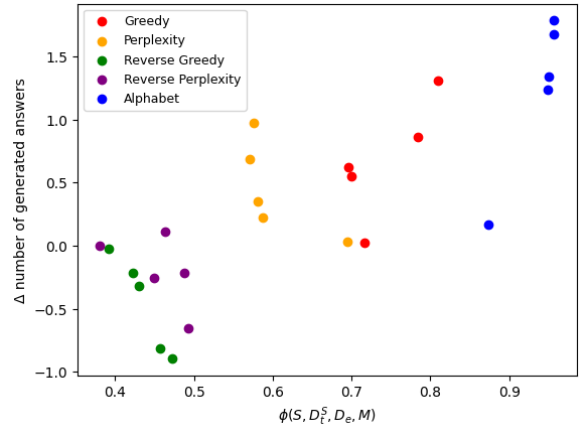


Figure 4: $\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M})$ vs. the number of generated answers across three datasets, where we use Llama2 for \mathcal{M} . Instead of the raw number of answer set, we report the size difference compared to the answer set generated from random ordering. As ϕ increases, which signifies how faithfully LM follows the ordering strategy in in-context examples, the model generates more answers.

$N_i = \sum_{j=1}^{m-1} \mathbb{1}(f(a_{ij}) > f(a_{i(j+1)}))$ represents the number of pairs that violates the ordering. Then, we compute micro average over \mathcal{D}_e .

$$\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M}) = \frac{100 \cdot \sum_{i \in \mathcal{D}_e} P_i}{\sum_{i \in \mathcal{D}_e} (P_i + N_i)}$$

Results Table 3 presents the results for Llama2 model, and we provide the results for OPT model in Table 8 in the appendix. For each $\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M})$, we also report $\phi(S, \mathcal{D}_t^{S_{\text{random}}}, \mathcal{D}_e, \mathcal{M})$ as a control. We found that in every cell (except for one cell in Table 8), the first number is higher than the second number, suggesting that the model follows the answer ordering pattern presented in the in-context examples. We found this is particularly true for ALPHABET ordering, which is probably the easiest pattern to learn.

<i>QAMPARI</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	26.3 / 25.2	11.7 / 10.9	13.8 / 12.9	25.3 / 22.4
GREEDY	26.4 / 25.7	12.2 / 11.9*	14.2 / 14.0*	25.6 / 22.6
PERPLEXITY	26.7 / 26.4*	12.4* / 11.6*	14.6* / 13.9*	25.8 / 22.9
REVERSE GREEDY	26.5 / 25.8	11.6 / 10.1*	13.9 / 12.4	25.1 / 21.8
REVERSE PERPLEXITY	27.0 / 26.7*	11.7 / 11.0	14.0 / 13.3	25.2 / 22.5
ALPHABET	24.5* / 23.5*	12.7* / 11.8*	14.3 / 13.6	24.7 / 22.6
<i>QUEST</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	23.9 / 24.8	17.9 / 19.7	18.3 / 19.9	27.2 / 27.8
GREEDY	23.8 / 24.8	19.6* / 20.8*	19.5* / 20.6*	28.6* / 28.4*
PERPLEXITY	24.3 / 24.8	19.3* / 20.8*	19.4 / 20.6*	28.0 / 28.4*
REVERSE GREEDY	22.9 / 24.5	17.0 / 18.4*	17.4 / 18.8*	26.3 / 26.5*
REVERSE PERPLEXITY	23.7 / 24.5	17.3 / 19.4	17.7 / 19.4	26.4 / 27.1*
ALPHABET	20.5* / 23.8*	17.6 / 20.4*	17.0 / 20.0	25.0* / 27.0*

Table 4: QA performance for answer ordering strategies on Llama2 (13B) model. P_{EM} and R_{EM} are precision and recall for calculating $F1_{EM}$. We present development set performance and then test set performance in each cell. Blue color indicates improved performance compared to Random and red indicates the opposite. We put * on scores that are significantly different from that of Random ordering.

We further observe that the model is decoding answers such that it will present **confident** answer first (following the orders of GREEDY and PERPLEXITY), even when answers in in-context example is randomly ordered. Even after introducing consistent ordering (presenting less confident answer first), the model shows propensity to present confident answer first (values for REVERSE GREEDY and REVERSE PERPLEXITY are below chance (50) consistently).

5.2 Does ordering impact the number of generated answers?

Unlike in simpler QA tasks where there is exactly one gold answer, models have to decide how many answers to generate. Would consistent ordering of answers allow the model to generate more answers?

We report the number of generated answers for each ordering strategy for Llama model in Figure 4. We find that generation order impacts the number of generated answer, with ALPHABET ordering substantially increasing the number of generated answers the most. The results further suggest that an ordering pattern that is easier for the model to learn can prompt LM to generate more answers. We report the results for OPT model in Figure 8 which shows the same trends.⁵

5.3 Does the ordering impact the QA performance?

Lastly, we examine the end task (QA) performance of different answer ordering strategies. Table 4

⁵We did not measure it for GPT-3.5 as it is costly.

presents the results on QAMPARI and QUEST datasets on Llama2 model. Overall, we see that answer ordering does not bring large impact in final performance, but notice consistent patterns. Presenting more confident answers first (GREEDY and PERPLEXITY) yielded better results than their REVERSE counterparts. GREEDY and PERPLEXITY show gains mostly in recall, leading to increase in both $F1_{EM}$ and $F1_{F1}$. Arbitrary, yet consistent ordering such as ALPHABET does not improve model performance, sometimes rather leading to lower performance. The trend holds for AmbigQA (results presented in Table 7 in the appendix) though not statistically significant. This might be caused by smaller average answer set size compared to that of other datasets (2-3 vs. 10+ answers). We suggest ordering ‘known’ answer first in in-context examples to improve model performance.

For OPT model (result can be found in Table 9 in the appendix), we observe GREEDY and PERPLEXITY show improved performance through gains in recall for QUEST dataset but the results are mostly random on other datasets. We plot the perplexity of individual answer in train examples with respect to two models in Figure 5. Overall, we find that Llama2 contains more factual knowledge than OPT, resulting in higher end task performance. Two models exhibit similar knowledge for QUEST as they strongly correlate, however OPT shows a wider range of perplexities on other datasets, especially for answers that have low perplexity on Llama2. We hypothesis carefully ordering between answers will bring significant changes in end task

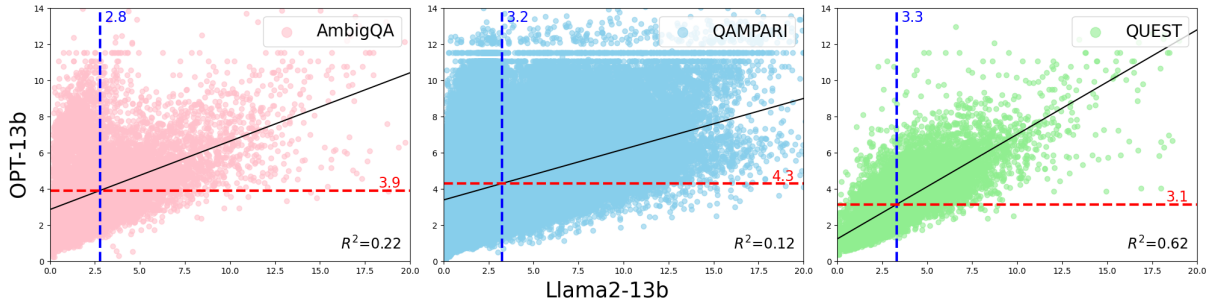


Figure 5: Plots of log answer perplexities from Llama2-13b (x-axis) and OPT-13b (y-axis). Horizontal and vertical lines indicate the mean value of log perplexities with respect to each LM. In all datasets, Llama2 outperforms OPT in its parametric knowledge, and the answers mostly report higher perplexity with OPT compared to Llama2.

performance **only** when model exhibits sufficient parametric knowledge of subset of answers. When the model is not familiar enough with the gold answers in in-context examples, knowledge-aware answer ordering might have limited effectiveness.

5.4 Transfer to other base LMs

So far we have measured the parametric knowledge on an language model and then use the same model for in-context prompting. In this section, we experiment using in-context example set constructed with parametric knowledge of one language model (Llama2), see how it impacts the generation of another language model (GPT-3.5). While different LMs have different pre-training data, the relative parametric knowledge might be similar for different LMs (e.g., famous entity to one LM remains famous for another LM). This also allows us to experiment with propriety black-box LM API easily, whose prediction probability is not always available. We observe similar patterns as in the original experiments (GPT 3.5 results in Table 10 in the appendix), but the effect size is much smaller and not significant, potentially because of the difference in parametric knowledge between two models.

6 Related Work

Analysis on In-context Learning Many prior works investigate factors that determine the performance of in-context learning (Brown et al., 2020), such as the composition of the pre-training dataset (Xie et al., 2022), size of language model (Wei et al., 2022a), number of pre-training tokens (Touvron et al., 2023), and specific fine-tuning strategy employed (Wei et al., 2021). More closely related to ours, one line of work particularly focuses on factors related to the in-context examples, including the choice of verbalizer and templates (Min et al., 2022), order of examples (Lu et al., 2022;

Pezeshkpour and Hruschka, 2023), and the choice of in-context examples (Liu et al., 2021; Rubin et al., 2021; Agrawal et al., 2022; Ye et al., 2023). While past work is mainly centered around classification tasks, our work studies the task of multi-answer QA, with a focus on how LM’s parametric knowledge on in-context examples impact the performance. In particular, our findings suggests that answers with lower perplexity lead to more accurate answer, which is congruent with recent work that shows using lower perplexity prompts improves model perplexity in general (Ye and Durrett, 2023; Iyer et al., 2023; Gonen et al., 2022).

Multi-answer QA Real-world questions could naturally have multiple answers when a question is ambiguous (Min et al., 2020; Stelmakh et al., 2022), when a question is evaluated under different temporal or geographical contexts (Zhang and Choi, 2021), or when a question expects a set of answers (Amouyal et al., 2022; Malaviya et al., 2023). While most prior work tackles multi-answer QA in the open-book setting by retrieving from external corpus (Shao and Huang, 2022; Sun et al., 2023), we study the problem in the close-book setting, which prompts LLMs to generate the answers based on their parametric knowledge.

7 Conclusion

We present comprehensive studies on knowledge-aware prompt design for multi-answer QA tasks. Our findings underscore the benefits of having in-context examples that the language model is familiar with. First, the HALFKNOWN set aids the model in effectively accessing its parametric knowledge. Second, employing knowledge-aware ordering of presenting answers in descending order of the model’s knowledge enhances the overall process of answer generation.

559
560
561
562
563
564
565
566
567
568

569

570
571
572
573
574

575
576
577
578
579

580
581
582
583
584
585

586
587
588
589
590
591

592
593
594

595
596
597
598

599
600
601
602

603
604
605
606

607
608
609
610

Limitations

Our study mainly focuses on multi-answer QA datasets. The analysis can be extended to a wide range of tasks that requires different types of reasoning ability. Also, we find that the end task performance gets less impacted when random in-context examples are used (Appendix F). Further studies can be conducted with diverse in-context example retrieval methods as well as cover multiple languages.

References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. In-context examples selection for machine translation. In *Annual Meeting of the Association for Computational Linguistics*.

Samuel Joseph Amouyal, Ohad Rubin, Ori Yoran, Tomer Wolfson, Jonathan Herzig, and Jonathan Berant. 2022. Qampari: An open-domain question answering benchmark for questions with many answers from multiple paragraphs. *ArXiv, abs/2205.12665*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *ArXiv, abs/2110.14168*.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *ArXiv, abs/2104.08821*.

Hila Gonen, Srini Iyer, Terra Blevins, Noah A. Smith, and Luke Zettlemoyer. 2022. Demystifying prompts in language models via perplexity estimation. *ArXiv, abs/2212.04037*.

Dan Iter, Reid Pryzant, Ruochen Xu, Shuo Wang, Yang Liu, Yichong Xu, and Chenguang Zhu. 2023. In-context demonstration selection with cross entropy difference. *ArXiv, abs/2305.14726*.

Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *ArXiv, abs/1705.03551*.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alben, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark

for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466. 611
612
613

Itay Levy, Ben Bogin, and Jonathan Berant. 2022. Diverse demonstrations improve in-context compositional generalization. *ArXiv, abs/2212.06800*. 614
615
616

Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*. 617
618
619
620

Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 621
622
623
624
625
626

Aman Madaan, Dheeraj Rajagopal, Niket Tandon, Yiming Yang, and Antoine Bosselut. 2022. Conditional set generation using seq2seq models. In *Conference on Empirical Methods in Natural Language Processing*. 627
628
629
630
631

Chaitanya Malaviya, Peter Shaw, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2023. Quest: A retrieval dataset of entity-seeking queries with implicit set operations. *arXiv preprint arXiv:2305.11694*. 632
633
634
635
636

Sewon Min, Xixi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*. 637
638
639
640
641

Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. Ambigqa: Answering ambiguous open-domain questions. *arXiv preprint arXiv:2004.10645*. 642
643
644
645

Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions. *ArXiv, abs/2308.11483*. 646
647
648
649

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*. 650
651
652

Zhihong Shao and Minlie Huang. 2022. Answering open-domain multi-answer questions via a recall-then-verify framework. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1825–1838, Dublin, Ireland. Association for Computational Linguistics. 653
654
655
656
657
658
659

Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. ASQA: Factoid questions meet long-form answers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8273–8288, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 660
661
662
663
664
665
666

667 Weiwei Sun, Hengyi Cai, Hongshen Chen, Pengjie Ren,
668 Zhumin Chen, Maarten de Rijke, and Zhaochun Ren.
669 2023. Answering ambiguous questions via iterative
670 prompting. In *Proceedings of the 61st Annual Meet-*
671 *ing of the Association for Computational Linguistics*
672 *(Volume 1: Long Papers)*, pages 7669–7683, Toronto,
673 Canada. Association for Computational Linguistics.

674 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
675 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
676 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti
677 Bhosale, et al. 2023. Llama 2: Open founda-
678 tion and fine-tuned chat models. *arXiv preprint*
679 *arXiv:2307.09288*.

680 Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin
681 Guu, Adams Wei Yu, Brian Lester, Nan Du, An-
682 drew M Dai, and Quoc V Le. 2021. Finetuned lan-
683 guage models are zero-shot learners. *arXiv preprint*
684 *arXiv:2109.01652*.

685 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,
686 Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
687 Maarten Bosma, Denny Zhou, Donald Metzler, et al.
688 2022a. Emergent abilities of large language models.
689 *arXiv preprint arXiv:2206.07682*.

690 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
691 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
692 et al. 2022b. Chain-of-thought prompting elicits rea-
693 soning in large language models. *Advances in Neural*
694 *Information Processing Systems*, 35:24824–24837.

695 Sang Michael Xie, Aditi Raghunathan, Percy Liang,
696 and Tengyu Ma. 2022. An explanation of in-context
697 learning as implicit bayesian inference. In *Internat-*
698 *ional Conference on Learning Representations*.

699 Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu,
700 and Lingpeng Kong. 2023. Compositional ex-
701 emplars for in-context learning. *arXiv preprint*
702 *arXiv:2302.05698*.

703 Xi Ye and Greg Durrett. 2023. Explanation selection
704 using unlabeled data for in-context learning. In *Pro-*
705 *ceedings of EMNLP*.

706 Michael Zhang and Eunsol Choi. 2021. SituatedQA: In-
707 corporating extra-linguistic contexts into QA. In *Pro-*
708 *ceedings of the 2021 Conference on Empirical Meth-*
709 *ods in Natural Language Processing*, pages 7371–
710 7387, Online and Punta Cana, Dominican Republic.
711 Association for Computational Linguistics.

712 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
713 Artetxe, Moya Chen, Shuohui Chen, Christopher
714 Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin,
715 Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shus-
716 ter, Daniel Simig, Punit Singh Koura, Anjali Srid-
717 har, Tianlu Wang, and Luke Zettlemoyer. 2022.
718 Opt: Open pre-trained transformer language mod-
719 els. *ArXiv*, abs/2205.01068.

720 Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and
721 Sameer Singh. 2021. Calibrate before use: Improv-
722 ing few-shot performance of language models. In
723 *International Conference on Machine Learning*.

A Dataset Statistics

We report the dataset statistics in Table 5.

B Similarity of In-Context Examples

We calculate the similarity score of two in-context examples using SimCSE embeddings of each query. Figure 6 illustrates the similarity distributions across three datasets.

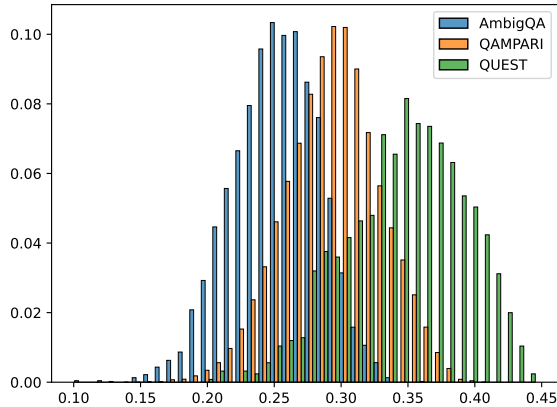


Figure 6: Similarity distributions among in-context example candidates. The x-axis denotes embedding similarity (with SimCSE (Gao et al., 2021) encoder) and the y-axis indicates the percentage of each bin. The median value for each dataset is 0.254, 0.295, 0.350.

C Experimental Details

C.1 Resources

All experiments are conducted on NVIDIA A40 GPU. A single evaluation for AmbigQA and QUEST (development split) took around 20 minutes. QAMPARI (development and test split) took around 1 hours. QUEST (test split) took around 2 hours, due to its largest size.

C.2 Statistical Testing

We conduct paired bootstrap tests with 10000 bootstrap samples throughout our experiments (Section 2.2). Since we have multiple (two or four) in-context example sets for experiments in Section 3, we randomly sample one in-context example set of each class (UNKNOWN, HALFKNOWN, KNOWN, and RANDOM) and conduct testing.

D In-Context Example Set Study

In Table 6, we present the results from Section 3.1 for QAMPARI_{test} and QUEST_{dev} on Llama2.

E Answer Ordering Strategies

E.1 Single Answer Study

We examine the effectiveness of answer ordering strategies discussed at Section 4. We provide only one answer at the forefront of each ordered answers in in-context examples. Since an answer from GREEDY and PERPLEXITY is ‘known’ to the model, they may serve as an upper bound of ‘known’ answer, while REVERSE GREEDY and REVERSE PERPLEXITY may serve as a lower bound. RANDOM exists somewhere between these. The disparities among these are clear, as shown in Figure 7. The results suggest that the model is able to differentiate ordering strategies.

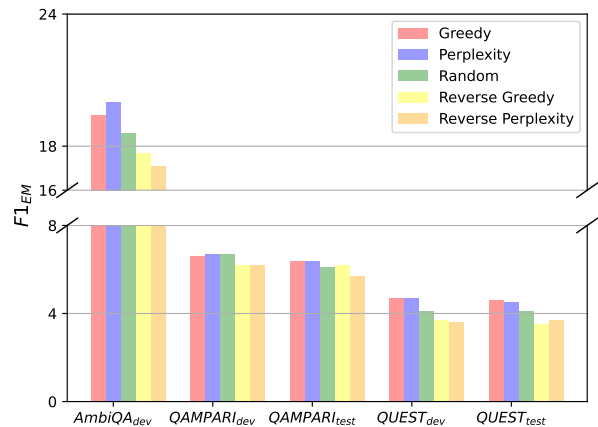


Figure 7: Answer-level Exact Match ($F1_{EM}$) score for demonstrating only one frontmost answer of an ordering methodology on Llama2 model.

E.2 AmbigQA results

We present the performance of answer ordering strategies on AmbigQA dataset in Table 7.

E.3 Results on OPT 13B model

We present the results of experiments in Section 5 with OPT 13B model. With respect to following the ordering strategy of in-context examples (Section 5.1, 5.2), we find that the results hold for OPT LLM model as well (Table 8). However, the end task performance results are somewhat mixed (Table 9, Figure 8). We observe consistent results of end task performance on QUEST dataset but the results are mostly random on AmbigQA and QAMPARI dataset.

E.4 Results on GPT-3.5 model

GPT-3.5-turbo model tends to generate lengthy and chatty outputs such as “There is not enough infor-

	AmbigQA		QAMPARI			QUEST		
	Train	Dev.	Train	Dev.	Test	Train	Dev.	Test
# Examples	4,615	1,048	50,372	1,000	1,000	1,251	316	1,669
Avg. # of answers	2.8	3.1	14.0	13.2	13.1	10.9	10.7	10.7
Query length	46.9	46.7	67.8	57.7	55.8	54.0	52.2	53.3
Answer length	15.9	14.5	14.4	17.3	16.6	17.2	16.7	17.0
Answer sequence length	45.2	45.4	200.9	228.5	217.6	187.0	179.0	182.4
# Unique answers	10,684	2,999	455,469	12,462	12,464	10,160	3,050	12,367

Table 5: Dataset statistics. Lengths of query, answer, and answer sequence are measured by the length of each string. # Unique answers counts unique answers within each split. Duplicated questions are removed from training sets.

	QAMPARI _{test}		QUEST _{dev}	
	$F1_{EM}$	$F1_{F1}$	$F1_{EM}$	$F1_{F1}$
Random	10.0	19.3	4.0	12.1
Unknown	10.6	20.2*	4.4*	13.2*
HalfKnown	11.2*	20.9*	4.9*	13.1*
Known	9.9	18.6	4.3*	12.8*

Table 6: Results comparing known example and unknown example. We put * on scores that are significantly different from that of Random in-context examples set, and bold the highest performing set for each metric.

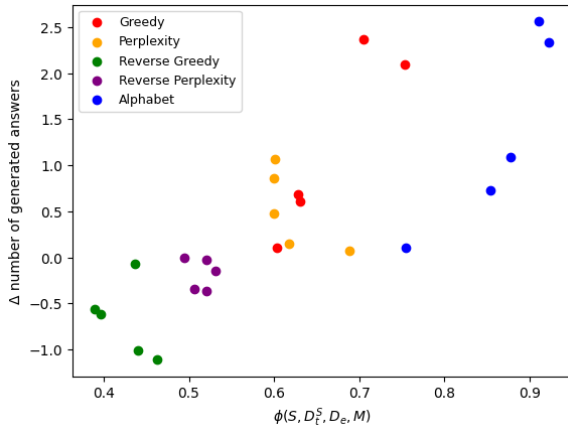


Figure 8: $\phi(S, \mathcal{D}_t^S, \mathcal{D}_e, \mathcal{M})$ vs. the number of generated answers across three datasets, where we use OPT (13B) model for \mathcal{M} .

mation given to answer this question". Therefore we add a short instruction as following: "Follow the answers pattern".

Table 10 shows the results for experiment in Section 5.4. We do not experiment on GREEDY and REVERSE GREEDY because we do not think that a greedy ordering will be effectively transferred between different LMs.

F Random Examples

Prior works have highlighted the importance of relevant in-context examples, such as those based on similarity (Liu et al., 2021) and diversity (Levy

et al., 2022). Yet, many studies do not do example specific retrieval and use random examples for its simplicity. Throughout our experiments (except for Section 3.1 which constructs universal in-context set for all examples in evaluation dataset), we retrieved similar in-context examples for each evaluation example. How would our results hold if we use randomly select in-context examples?

First, with randomly retrieved in-context examples, models still learn to follow the answer ordering strategy shown in in-context examples but substantially less than when using similar in-context examples (Table 11). Second, we find that the number of generated answer is affected similarly, with using ALPHABET ordering leads to the highest number of generated answers. However, we see invariant performances on end tasks (Table 12). Carefully constructing relevant in-context examples is more meaningful than doing it for random in-context examples. This suggests that if you do not have large enough training examples to recover semantically relevant in-context examples, careful construction of prompt might not yield changes in end task performance.

G Prompts

Throughout Table 13 to Table 16, we present the prompts used in our experiments.

<i>AmbigQA</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	27.1	17.9	20.0	31.3
GREEDY	27.2	18.5	20.5	31.7
PERPLEXITY	27.4	18.4	20.5	31.8
REVERSE GREEDY	27.1	17.8	20.1	31.5
REVERSE PERPLEXITY	27.3	17.9	20.2	31.8
ALPHABET	26.7	18.2	20.3	31.2

Table 7: QA performance on AmbigQA dataset on Llama2 model. The table is formatted the same as Table 4.

		<i>S</i>				
		GREEDY	REVERSE GREEDY	PERPLEXITY	REVERSE PERPLEXITY	ALPHABET
\mathcal{D}_e	AmbigQA _{dev}	60.3 / 58.3	43.7 / 42.2	68.8 / 58.1	49.5 / 41.9	75.5 / 50.5
	QAMPARI _{dev}	62.8 / 52.1	39.0 / 39.6	60.0 / 55.1	52.1 / 44.9	87.8 / 52.0
	QAMPARI _{test}	63.1 / 52.4	39.7 / 39.1	61.8 / 56.7	52.1 / 39.1	85.4 / 47.3
	QUEST _{dev}	70.5 / 49.1	44.0 / 42.5	60.0 / 57.1	53.1 / 42.9	91.1 / 67.6
	QUEST _{test}	75.3 / 57.5	46.3 / 45.5	60.1 / 54.0	50.6 / 46.0	92.3 / 51.6
Average		66.4	42.5	62.1	51.5	86.4

Table 8: Percentage of generated answer ordering matching in-context examples answer ordering, where we use OPT (13B) model for \mathcal{M} . The table is formatted the same as Table 3.

<i>AmbigQA</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	13.1	10.3	10.7	19.4
GREEDY	13.1	10.3	10.7	19.5
PERPLEXITY	12.9	10.0	10.5	19.2
REVERSE GREEDY	12.9	9.9	10.5	19.1
REVERSE PERPLEXITY	13.2	10.7	11.0	19.3
ALPHABET	13.5	10.6	11.0	19.3
<i>QAMPARI</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	14.2 / 15.5	7.5 / 7.2	8.1 / 8.2	18.6 / 17.1
GREEDY	14.0 / 14.9	7.5 / 7.6	7.9 / 8.4	18.6 / 17.8
PERPLEXITY	14.7 / 15.6	7.8 / 7.7	8.3 / 8.5	19.0 / 17.6
REVERSE GREEDY	14.5 / 15.4	6.9* / 6.7	7.6 / 7.9	18.0* / 16.7
REVERSE PERPLEXITY	15.6* / 15.9	7.6 / 7.2	8.4 / 8.3	18.8 / 16.9
ALPHABET	14.4 / 15.0	8.1* / 7.9*	8.5 / 8.9*	18.7 / 17.4
<i>QUEST</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	14.6 / 18.4	11.6 / 16.1	12.0 / 15.6	21.3 / 23.8
GREEDY	15.7 / 18.6	16.6* / 18.0*	14.9* / 17.0*	23.7* / 25.2*
PERPLEXITY	16.1 / 18.3	14.8* / 17.0*	13.9* / 16.2	22.6 / 24.5
REVERSE GREEDY	14.5 / 17.4*	10.7 / 13.8*	10.2 / 13.8*	19.6 / 22.1*
REVERSE PERPLEXITY	15.0 / 17.9	14.3* / 15.4	13.2 / 15.1	22.4 / 23.5
ALPHABET	16.3* / 17.6*	15.9* / 17.3*	14.7* / 16.3*	23.0* / 24.1*

Table 9: QA performance for answer ordering strategies with OPT (13B) model. The table is formatted the same as Table 4.

<i>AmbigQA</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	28.2	22.1	23.1	35.7
PERPLEXITY	28.8	23.1*	23.9	36.5
REVERSE PERPLEXITY	29.0	22.3	23.5	35.3
ALPHABET	28.4	22.5	23.5	35.8
<i>QAMPARI</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	23.4 / 23.2	18.7 / 18.5	18.4 / 18.4	30.1 / 28.4
PERPLEXITY	23.9 / 22.9	19.5 / 19.1	18.9 / 18.5	30.4 / 29.1
REVERSE PERPLEXITY	23.2 / 23.1	18.2 / 18.5	18.2 / 18.3	30.2 / 28.5
ALPHABET	23.4 / 23.0	17.3* / 17.8	17.8 / 18.0	29.0* / 27.5
<i>QUEST</i>	P_{EM}	R_{EM}	$F1_{EM}$	$F1_{F1}$
RANDOM	15.0 / 16.4	16.7 / 17.6	14.8 / 15.8	25.5 / 26.4
PERPLEXITY	16.6 / 17.0	17.7 / 18.6*	15.9 / 16.5*	26.8 / 26.8
REVERSE PERPLEXITY	16.2 / 16.5	17.5 / 17.8	15.5 / 15.9	26.6 / 26.4
ALPHABET	15.5 / 17.0	16.2 / 17.6	14.9 / 16.2	24.9 / 25.5*

Table 10: QA performance for answer ordering strategies with GPT-3.5 model. We use the particular answer ordering from Llama2 and transfer to GPT-3.5 model. The table is formatted the same as Table 4 except we do not experiment on GREEDY and REVERSE GREEDY.

	<i>S</i>				
	GREEDY	REVERSE GREEDY	PERPLEXITY	REVERSE PERPLEXITY	ALPHABET
\mathcal{D}_e AmbigQA _{dev}	69.6 / 68.9	33.7 / 32.8	70.2 / 70.5	68.9 / 29.5	83.6 / 62.5
QAMPARI _{dev}	63.2 / 59.8	40.7 / 40.3	57.0 / 57.3	57.3 / 42.7	92.6 / 65.9
QAMPARI _{test}	61.2 / 61.4	43.5 / 43.2	57.5 / 56.4	57.5 / 43.6	92.7 / 60.7
QUEST _{dev}	55.4 / 52.6	39.3 / 40.2	59.1 / 57.4	57.1 / 42.6	88.5 / 59.7
QUEST _{test}	56.8 / 54.0	38.9 / 40.1	56.9 / 56.4	56.4 / 43.6	86.7 / 60.9
Average	61.2	39.2	60.1	59.4	88.8

Table 11: Percentage of generated answer ordering matching in-context examples answer ordering, where we employ **random** in-context examples instead of most similar examples. The table is formatted the same as Table 3.

	<i>AmbigQA_{dev}</i>			<i>QAMPARI_{dev}</i>			<i>QAMPARI_{test}</i>		
	$F1_{EM}$	$F1_{F1}$	# ans	$F1_{EM}$	$F1_{F1}$	# ans	$F1_{EM}$	$F1_{F1}$	# ans
RANDOM	17.8	28.7	2.07	9.8	20.2	3.77	10.0	19.1	3.74
GREEDY	17.4	27.8	2.12	9.6	19.9	4.42	9.3	17.7	4.43
PERPLEXITY	17.9	28.3	2.11	9.7	20.0	3.99	9.7	18.6	4.03
REVERSE GREEDY	17.6	28.3	2.08	9.8	20.4	3.82	9.6	18.5	3.61
REVERSE PERPLEXITY	17.9	28.4	2.11	9.3	19.7	3.83	9.6	18.4	3.81
ALPHABET	17.9	28.5	2.22	9.8	19.8	5.48	9.6	17.5	5.41

	<i>QUEST_{dev}</i>			<i>QUEST_{test}</i>		
	$F1_{EM}$	$F1_{F1}$	# ans	$F1_{EM}$	$F1_{F1}$	# ans
RANDOM	4.4	12.9	3.42	3.5	11.2	3.41
GREEDY	4.7	12.5	4.51	3.4	10.9	4.49
PERPLEXITY	4.7	13.0	3.60	3.4	11.1	3.62
REVERSE GREEDY	4.0	12.5	3.51	3.3	11.1	3.11
REVERSE PERPLEXITY	4.6	12.6	3.09	3.6	11.4	3.28
ALPHABET	4.5	11.2	5.84	3.0	9.4	5.99

Table 12: QA performance for answer ordering strategies with random in-context examples. We bold the highest performing set for each metric.

Question	Who is the current chairman of african union commission?
Gold Answers	Jean Ping, Moussa Faki, Nkosazana Clarice Dlamini-Zuma
Prompt	Question: Who is the chairman of the federal reserve?\nAnswers: Alan Greenspan Ben Bernanke Janet Yellen\n\nQuestion: Who is the president of south africa now?\nAnswers: Thabo Mvuyelwa Mbeki Kgalema Petrus Motlanthe JZ\n\nQuestion: Who is the present chairperson of national human rights commission in india?\nAnswers: Justice K. G. Balakrishnan H. L. Dattu Cyriac Joseph\n\nQuestion: Who appoints the chairman of the finance commission?\nAnswers: the President Pranab Mukherjee Ram Nath Kovind Pratibha Devisingh Patil\n\nQuestion: Who is the chairman of national commission for woman of india?\nAnswers: Lalitha Kumaramangalam Mamta Sharma Girija Vyas\n\nQuestion: Who is the current chairman of african union commission?\nAnswers:
Output	Jean Ping Nkosazana Dlamini-Zuma Moussa Faki Mahamat\n

Table 13: Prompt example of AmbigQA

Question	What movies did Scott Z. Burns screenwrite?
Gold Answers	Contagion, No Time to Die, Pu-239, Side Effects, The Bourne Ultimatum, The Informant!, The Laundromat, The Mercy, The Report
Prompt	Question: Which film has Edward Burns as a member of its cast and had Edward Burns as screenwriter?\nAnswers: Ash Wednesday Purple Violets She’s the One Sidewalks of New York The Brothers McMullen The Groomsman Newlyweds Looking for Kitty No Looking Back\n\nQuestion: Scott Ziehl was a director for what genre of film?\nAnswers: action film crime thriller horror film monster film drama\n\nQuestion: What are the publication dates of film that had Scott Z. Burns as screenwriter?\nAnswers: 2007 2006 2009 2013 2014 2018 2019 2020\n\nQuestion: Who are the cast members of film that had Scott Z. Burns as screenwriter?\nAnswers: Christoph Waltz Patton Oswalt Marsha Stephanie Blake Jocko Sims David Costabile Dan Fredenburgh Ted Levine Ken Stott Steven Berkoff\n\nQuestion: What are the genres of film that had Scott Z. Burns as screenwriter?\nAnswers: action film thriller film comedy film spy film drama\n\nQuestion: What movies did Scott Z. Burns screenwrite?\nAnswers:
Output	Contagion Side Effects The Bourne Ultimatum The Bourne Supremacy The Bourne Identity\n

Table 14: Prompt example of QAMPARI

Question	Mary Stewart novels
Gold Answers	A Walk in Wolf Wood, Airs Above the Ground, My Brother Michael, The Crystal Cave, The Hollow Hills, The Ivy Tree, The Last Enchantment, The Prince and the Pilgrim, The Wicked Day, The Wind Off the Small Isles, This Rough Magic, Thornyhold, Touch Not the Cat, Wildfire at Midnight
Prompt	Question: 1740s fantasy novels, 1747 books, or novels by Denis Diderot\nAnswers: Niels Klim’s Underground Travels The Sofa: A Moral Tale Jacques the Fatalist Rameau’s Nephew La Religieuse The Art of Cookery Made Plain and Easy Letters from a Peruvian Woman Madame de La Carli\n\nQuestion: Sarah Weeks novels\nAnswers: So B. It Jumping the Scratch\n\nQuestion: 1858 British novels, 1872 fantasy novels, or Novels by George MacDonald\nAnswers: At the Back of the North Wind Alec Forbes of Howglen The Princess and the Goblin David Elginbrod The Day Boy and the Night Girl Randal Bannerman’s Boyhood The Princess and Curdie Through the Looking-Glass Phantastes My Lady Ludlow Sir Gibbie Mary Marston Doctor Thorne Lilith Poor Dionis\n\nQuestion: Mark Clapham novels\nAnswers: The Taking of Planet 5 Twilight of the Gods The Book of the War Beige Planet Mars Hope\n\nQuestion: Marlon James Novels\nAnswers: A Brief History of Seven Killings John Crow’s Devil The Book of Night Women\n\nQuestion: Mary Stewart novels\nAnswers:
Output	The Crystal Cave The Hollow Hills The Last Enchantment The Wicked Day\n

Table 15: Prompt example of QUEST

Question	Every day, Wendi feeds each of her chickens three cups of mixed chicken feed, containing seeds, mealworms and vegetables to help keep them healthy. She gives the chickens their feed in three separate meals. In the morning, she gives her flock of chickens 15 cups of feed. In the afternoon, she gives her chickens another 25 cups of feed. How many cups of feed does she need to give her chickens in the final meal of the day if the size of Wendi's flock is 20 chickens?
Logical reasoning and Answer	If each chicken eats 3 cups of feed per day, then for 20 chickens they would need $3 \times 20 = 60$ cups of feed per day. If she feeds the flock 15 cups of feed in the morning, and 25 cups in the afternoon, then the final meal would require $60 - 15 - 25 = 20$ cups of chicken feed. ##### 20
Prompt	Question: Mabel lives 4500 steps directly east of Lake High school. Helen lives $\frac{3}{4}$ the number of steps that Mabel lives, directly west of the school. What's the total number of steps Mabel will walk to visit Helen so that they can do their assignments together? Answer: Helen lives $\frac{3}{4} \times 4500 = 3375$ steps directly west of Lake High. To reach Helen, Mabel would have to walk to $4500 + 3375 = 7875$ steps. ##### 7875 Question: Mark is 7 years older than Amy, who is 15. How old will Mark be in 5 years? Answer: Mark is 15 years + 7 years = 22 years old. In 5 years, he will be 22 years + 5 years = 27 years old. ##### 27 Question: Steve has 2 boxes of pencils with 12 pencils in each box. He gave Matt 3 more pencils than he gave to Lauren. If Steve gave 6 pencils to Lauren, how many pencils does he have left? Answer: Steve started with $2 \times 12 = 24$ pencils. He gave Matt $6 + 3 = 9$ pencils. After giving away the pencils, Steve will have $24 - 9 = 15$ pencils left. ##### 9 Question: Mandy researched 42 med schools. She applied to $\frac{1}{3}$ of the schools she researched and got into half of the schools where she applied. How many schools was Mandy accepted to? Answer: First find the number of schools Mandy applied to: $42 \text{ med schools} \div 3 = 14$ med schools Then divide that number by 2 to find the number of schools where she was accepted: $14 \text{ med schools} \div 2 = 7$ med schools ##### 7 Question: Rachel is stuffing envelopes. She has eight hours to complete the task, and there are 1,500 envelopes. In the first hour, Rachel stuffs 135 envelopes. The second hour she stuffs 141 envelopes. How many envelopes will Rachel need to stuff per hour to finish the job? Answer: Rachel has $1500 - 135 - 141 = 1224$ envelopes remaining to stuff. Rachel has $8 - 2 = 6$ hours left to finish the task. Rachel needs to stuff $1224 \text{ envelopes} \div 6 \text{ hours} = 204$ envelopes per hour. ##### 204 Question: Samson is going to another town which is 140 km away. He will use his car that uses ten liters of gasoline for a distance of 70 km. How many liters of gasoline will Samson need for a one-way trip? Answer: Samson will need $140 \text{ km} \div 70 \text{ km} = 2$ ten liters of gasoline for a one-way trip to a town. Therefore, he will need a total of $2 \times 10 = 20$ liters of gasoline. ##### 20 Question: Every day, Wendi feeds each of her chickens three cups of mixed chicken feed, containing seeds, mealworms and vegetables to help keep them healthy. She gives the chickens their feed in three separate meals. In the morning, she gives her flock of chickens 15 cups of feed. In the afternoon, she gives her chickens another 25 cups of feed. How many cups of feed does she need to give her chickens in the final meal of the day if the size of Wendi's flock is 20 chickens? Answer:
Output	Wendi gives her chickens 15 cups of feed in the morning and 25 cups of feed in the afternoon. She needs to give her chickens another 20 cups of feed in the final meal of the day. ##### 20

Table 16: Prompt example of GSM8K