

In-Context Learning with Iterative Demonstration Selection

Anonymous ACL submission

Abstract

Spurred by advancements in scale, large language models (LLMs) have demonstrated strong few-shot learning ability via in-context learning (ICL). However, the performance of ICL has been shown to be highly sensitive to the selection of few-shot demonstrations. Selecting the most suitable examples as context remains an ongoing challenge and an open problem. Existing literature has highlighted the importance of selecting examples that are diverse or semantically similar to the test sample while ignoring the fact that the optimal selection dimension, *i.e.*, diversity or similarity, is task-specific. Leveraging the merits of both dimensions, we propose Iterative Demonstration Selection (IDS). Using zero-shot chain-of-thought reasoning (Zero-shot-CoT), IDS iteratively selects examples that are diverse but still strongly correlated with the test sample as ICL demonstrations. Specifically, IDS applies Zero-shot-CoT to the test sample before demonstration selection. The output reasoning path is then used to choose demonstrations that are prepended to the test sample for inference. The generated answer is followed by its corresponding reasoning path for extracting a new set of demonstrations in the next iteration. After several iterations, IDS adopts majority voting to obtain the final result. Through extensive experiments on tasks including reasoning, question answering, topic classification, and sentiment analysis, we demonstrate that IDS can consistently outperform existing ICL demonstration selection methods.

1 Introduction

With the recent advancements in scaling up model parameters, large language models (LLMs) showcase promising results on a variety of few-shot tasks through in-context learning (ICL), where the model is expected to directly generate the output of the test sample without updating parameters. This is achieved by conditioning on a manually designed

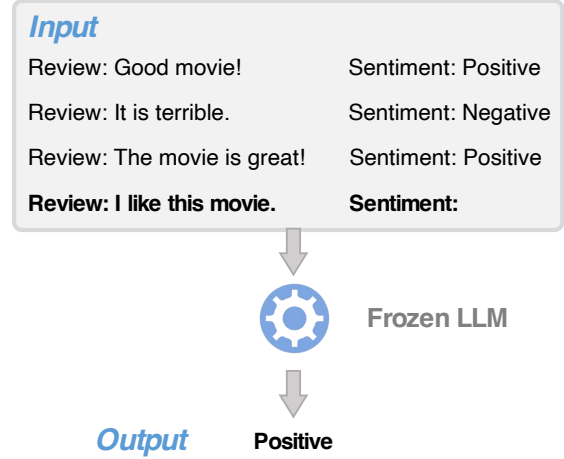


Figure 1: Illustration of in-context learning (ICL) on sentiment analysis. A frozen large language model directly generates the sentiment ‘Positive’ for the test sample ‘I like this movie.’ by taking the demonstrations and the test sample as input.

prompt consisting of an optional task description and a few demonstration examples (Brown et al., 2020). Fig. 1 shows an example describing how LLMs perform ICL on the sentiment analysis task. Given a few review-sentiment pairs as demonstrations, ICL combines them with the test sample as input, to the LLM for inference. The output, *i.e.*, ‘Positive’, is generated by the model autoregressively without any parameter updates.

Despite the effectiveness, the performance of ICL has been shown to be highly sensitive to the selection of demonstration examples (Zhao et al., 2021). Different sets of demonstrations can yield performance ranging from nearly random to comparable with state-of-the-art models (Gao et al., 2021; Lu et al., 2022). To alleviate the above issue, researchers in ICL have proposed a number of methods to select a set of examples as few-shot demonstrations (Rubin et al., 2022; Liu et al., 2022; Li and Qiu, 2023; Wang et al., 2023b; Li et al., 2023a; Ma et al., 2023; An et al., 2023b). Neverthe-

less, most of the existing approaches are only applicable to small language models as they typically require accessing model parameters or detailed output distributions which are usually not available for LLMs (Sun et al., 2022). Therefore, it is still a common practice to randomly select examples or select examples that are semantically similar to the test sample as demonstrations for LLMs, *i.e.*, considering diversity or similarity. While several approaches investigate the combination of similarity and diversity when prompting with explanations or exploring compositional generalization (Ye et al., 2022; An et al., 2023a), it is not yet clear how to determine and leverage the optimal dimension for different tasks in ICL.

Actually, the optimal dimension for selecting demonstration examples is task-specific. As we will show in §4, the diversity dimension is superior to the similarity dimension on CommonsenseQA and SST2 while the similarity dimension outperforms the diversity dimension on AGNews and BoolQ. Thus, it is unreasonable to claim that one dimension is consistently better than the other across different tasks.

To fully leverage the merits of both dimensions, we propose Iterative Demonstration Selection (IDS) for ICL (Fig. 2). IDS can iteratively select demonstration examples that are diverse but still have a strong correlation with the test sample through zero-shot chain-of-thought reasoning (Zero-shot-CoT) (Kojima et al., 2022). Specifically, Zero-shot-CoT, *e.g.*, “Let’s think step by step.”, is first applied to the test sample before selecting demonstrations to obtain a reasoning path. The training examples that are most semantically similar to the generated reasoning path are then selected as demonstrations. They are prepended to the test sample for inference. Note that IDS ensures that the generated answer is accompanied by the reasoning path through designed prompts. The new reasoning path is then used for extracting another set of demonstration examples by semantic similarity in the next iteration. After a few iterations, IDS adopts majority voting to obtain the final result. Empirical results on tasks spanning reasoning, question answering, topic classification, and sentiment analysis show that IDS can consistently outperform previous ICL demonstration selection baselines. In summary, our main contributions are:

- We consider both the diversity and similarity dimensions of ICL demonstration selection for

LLMs. We identify that the optimal dimension for selecting demonstrations is task-specific and propose Iterative Demonstration Selection (IDS) to fully leverage the merits of both dimensions.

- With extensive experiments and analysis, we demonstrate the effectiveness of IDS on a variety of tasks. Our code base is available at [<redacted>](#).

2 Related Work

This work mainly explores how to select few-shot in-context learning demonstrations for LLMs by leveraging Zero-shot-CoT. In light of this, we review four lines of research that form the basis of this work: few-shot learning, in-context learning basics, demonstration selection for in-context learning, and chain-of-thought reasoning.

2.1 Few-shot Learning

Few-shot learning aims to learn tasks with only a few labeled samples, which results in a big challenge, *i.e.*, over-fitting, for models as they typically require large amounts of data for training. Prior methods to address over-fitting mainly focused on augmenting the few-shot data (Gao et al., 2020; Qin and Joty, 2022), reducing the hypothesis space (Triantafillou et al., 2017; Hu et al., 2018), or optimizing the strategy for searching the best hypothesis (Ravi and Larochelle, 2017; Finn et al., 2017). More recently, LLMs have demonstrated strong few-shot learning ability through in-context learning without any parameter updates (Brown et al., 2020).

2.2 In-context Learning

Brown et al. (2020) first showed that a frozen GPT-3 model can achieve impressive results on a variety of few-shot NLP tasks through conditioning on manually designed prompts consisting of task descriptions and several demonstration examples. Since then many efforts have been made on in-context learning (ICL) (Dong et al., 2022). Chen et al. (2022); Min et al. (2022a); Wei et al. (2023a) demonstrated that the ICL ability of language models can be further improved through self-supervised or supervised training. Some analytical studies attempted to understand what factors affect ICL performance (Zhao et al., 2021; Shin et al., 2022; Wei et al., 2022a; Min et al., 2022b; Yoo et al., 2022; Wei et al., 2023b) and why ICL works (Xie et al., 2022; Olsson et al., 2022; Li et al., 2023b; Pan et al.,

2023; Dai et al., 2023). Other ongoing research on ICL has also explored (i) demonstration designing, including demonstration selection (Liu et al., 2022; Rubin et al., 2022; Wang et al., 2023b), demonstration ordering (Lu et al., 2022), and demonstration formatting (Wei et al., 2022b; Wang et al., 2022c; Zhou et al., 2023; Zhang et al., 2023a), (ii) applications of ICL (Ding et al., 2022; Meade et al., 2023; Zheng et al., 2023), and (iii) ICL beyond text (Wang et al., 2023c; Huang et al., 2023; Zhu et al., 2023; Wang et al., 2023a).

2.3 Demonstration Selection for In-context Learning

The performance of ICL has been shown to be highly sensitive to the selection of demonstration examples (Zhao et al., 2021). Existing methods to solve this problem can be mainly divided into two categories. First, *unsupervised* methods rely on pre-defined metrics. Liu et al. (2022) proposed to select the closest neighbors as demonstrations. In contrast, Levy et al. (2022) selected diverse demonstrations to improve in-context compositional generalization. More recent studies have explored leveraging the output distributions of language models to select few-shot demonstrations (Wu et al., 2022; Nguyen and Wong, 2023; Li and Qiu, 2023). Second, *supervised* methods involve model training. Rubin et al. (2022); Ye et al. (2023); Li et al. (2023a); Luo et al. (2023) proposed to learn to retrieve demonstration examples. Wang et al. (2023b) posited LMs as implicit topic models to facilitate demonstration selection. In addition, some studies (Zhang et al., 2022; Scarlatos and Lan, 2023) attempted to select demonstrations based on reinforcement learning. However, most of the existing methods are not applicable to LLMs as model parameters or output distributions are typically not available for LLMs (Sun et al., 2022), which motivates us to propose our simple but effective approach (IDS).

2.4 Chain-of-Thought Reasoning

Chain-of-thought (CoT) reasoning induces LLMs to produce intermediate reasoning steps before generating the final answer (Wei et al., 2022b). Depending on whether there are manually designed demonstrations, current CoT reasoning methods mainly include Manual-CoT and Zero-shot-CoT. In Manual-CoT, human-labeled reasoning paths are used to perform CoT reasoning (Wei et al., 2022b; Zhou et al., 2022; Wang et al., 2022b;

	CommonsenseQA	BoolQ	AGNews	SST2
Similar-ICL-Consistency (Similarity)	76.0	85.0	90.0	94.0
Random-ICL-Voting (Diversity)	79.0	84.0	88.0	95.0

Table 1: Results of different methods on CommonsenseQA, BoolQ, AGNews and SST2. The optimal dimension for selecting ICL demonstrations is task-specific.

Li et al., 2022; Wang et al., 2022a). In contrast, LLMs leverage self-generated rationales for reasoning in Zero-shot-CoT (Kojima et al., 2022; Zelikman et al., 2022; Zhang et al., 2023a; Diao et al., 2023). The ongoing research on CoT reasoning has also explored (i) multimodal reasoning (Zhang et al., 2023b; Wu et al., 2023), (ii) distilling knowledge from LLMs (Ho et al., 2022; Fu et al., 2023), and (iii) iterative optimization (Shinn et al., 2023; Madaan et al., 2023; Paul et al., 2023).

3 Problem Formulation

Given the test set $\mathcal{D}_{\text{test}}$ and the training set $\mathcal{D}_{\text{train}}$, the goal of ICL demonstration selection is to find an optimal subset $\mathcal{S} = \{(x_1, y_1), \dots, (x_k, y_k)\}$ (k -shot) of $\mathcal{D}_{\text{train}}$ as demonstration examples for each test sample (\hat{x}_i, \hat{y}_i) to maximize the overall task performance on $\mathcal{D}_{\text{test}}$. More formally, the optimal selection method \tilde{h} is defined as:

$$\tilde{h} = \arg \max_{h \in \mathcal{H}} \sum_{i=1}^{|\mathcal{D}_{\text{test}}|} \delta_{\text{LLM}([h(\mathcal{D}_{\text{train}}, \hat{x}_i, \hat{y}_i), \hat{x}_i]), \hat{y}_i} \quad (1)$$

where \mathcal{H} is the hypothesis space for searching demonstration examples, $h(\mathcal{D}_{\text{train}}, \hat{x}_i, \hat{y}_i)$ refers to demonstrations selected for (\hat{x}_i, \hat{y}_i) using h , $[,]$ stands for concatenation, and $\delta_{a,b}$ is the Kronecker delta function: $\delta_{a,b} = 1$ if a equals b , otherwise $\delta_{a,b} = 0$. In this work, we aim to find the optimal method \tilde{h} by leveraging Zero-shot-CoT.

4 What Makes Good In-Context Demonstrations?

As demonstrated in previous work (Zhao et al., 2021), the overall task performance is highly sensitive to the selection method h . Different sets of demonstration examples can yield significantly different performance. For example, Zhang et al. (2022) show that the minimum and maximum ICL performance due to random sampling differs by $> 30\%$ on 4 classification tasks, which emphasizes the importance of selecting good demonstrations for LLMs.

Similar-ICL-Consistency	Random-ICL-Voting
Which choice is the correct answer to the question?	Which choice is the correct answer to the question?
Examples: Question: If you have cleaned off dust here it may be difficult to do your homework where? Answer Choices: (A) desktop (B) closet (C) most buildings (D) surface of earth (E) stove Answer: A Question: Where is dust likely to be under? Answer Choices: (A) closet (B) ground (C) windowsill (D) attic (E) carpet Answer: E Question: Where would you find a dustbin that is being used? Answer Choices: (A) utility closet (B) ground (C) cupboard (D) broom closet (E) kitchen Answer: E Question: Dust accumulates where? Answer Choices: (A) ceiling (B) library (C) surface of earth (D) <u>most buildings</u> (E) desktop Answer: D	Examples: Question: She had a busy schedule, she had to run errands and pick up the kids the second she did what? Answer Choices: (A) make time for (B) take money (C) go outdoors (D) leave work (E) field Answer: D Question: What is the worst outcome of an injury? Answer Choices: (A) cause death (B) cause bleeding (C) falling down (D) become infected (E) claim insurance Answer: A Question: Mom said that Sarah should stay in bed until she was able to go to school again. What did mom say to Sarah when she tried to get up? Answer Choices: (A) you're sick (B) were sick (C) more rest (D) rest more (E) get back under the covers Answer: A Question: John got a raise, but he lost rank. Overall, it was a good what? Answer Choices: (A) demotion (B) push down (C) go off strike (D) lower (E) go off strike Answer: A
The response should follow the format: Answer: {A, B, C, D or E} Here is the test data. Question: John wanted to clean all of the dust out of his place before settling down to watch his favorite shows. What might he hardest do dust? Answer Choices: (A) closet (B) under the bed (C) television (D) attic (E) <u>most buildings</u>	The response should follow the format: Answer: {A, B, C, D or E} Here is the test data. Question: John wanted to clean all of the dust out of his place before settling down to watch his favorite shows. What might he hardest do dust? Answer Choices: (A) closet (B) under the bed (C) television (D) attic (E) most buildings
Answer: E ✗	Answer: D ✓

Table 2: Examples of Similar-ICL-Consistency (first decoding path) and Random-ICL-Voting (first run) for constructing demonstration examples. The upper part is the input to LLMs, including few-shot demonstrations, and the lower part is the predicted answer. Similar-ICL-Consistency gives the wrong answer ‘most buildings’ which is actually the output of the final demonstration example, indicating that the decision process of the model is misled by this similar sample.

A natural question is: what makes good in-context demonstrations? For LLMs, it is still a common practice to select a subset \mathcal{S} consisting of examples that are diverse or semantically similar to the test sample as demonstrations, *i.e.*, considering the diversity or similarity of \mathcal{S} . To investigate whether one dimension is consistently better than the other one across different tasks, we conduct some pilot experiments on CommonsenseQA (Talmor et al., 2019), BoolQ (Clark et al., 2019), AGNews (Zhang et al., 2015) and SST2 (Socher et al., 2013). Specifically, we randomly sample 100 examples from the original test set for experiments and conduct 4-shot learning using GPT-3.5 (gpt-3.5-turbo).

Following Zhang et al. (2023a), we use SentenceBERT (Reimers and Gurevych, 2019) to encode all samples. For each test sample, the Similar-ICL method selects the top-4 similar training data based on cosine similarity while the Random-ICL method randomly samples 4 training examples as

few-shot demonstrations. Inspired by Wang et al. (2022b), we apply *self-consistency* with 3 decoding paths (temperature 0.7) to Similar-ICL (named **Similar-ICL-Consistency**) and run Random-ICL 3 times before majority voting (named **Random-ICL-Voting**) to improve the robustness.

The results of different methods on four datasets are reported in Table 1. We can observe that the diversity dimension outperforms the similarity dimension on CommonsenseQA and SST2 while the similarity dimension is superior to the diversity dimension on BoolQ and AGNews. Therefore, the optimal dimension for selecting demonstration examples is task-specific. Thus, it is unreasonable to claim that one dimension is consistently better than the other one in ICL demonstration selection.

Intuitively, semantically similar examples can help the model correctly answer the test query as they might share similar input-output patterns with the test sample which could unleash GPT-3.5’s power of text generation. To further under-

stand why the similarity dimension underperforms the diversity dimension on CommonsenseQA, we present a case study in Table 2. We can see that the answer of the final demonstration example extracted by Similar-ICL-Consistency, *i.e.*, ‘most buildings’ is also in the options list of the test sample, which misleads the decision process of the model, leading to a wrong answer. In addition, the selected demonstrations might not include enough important information as high similarity also results in redundancy.

Considering the strengths and weaknesses of both dimensions, we aim to design a method that can select demonstration examples that are diverse (minimizing misleading information) but still strongly correlated with the test sample, which is introduced in the next section.

5 Iterative Demonstration Selection

Based on the observations and considerations in §4, we introduce Iterative Demonstration Selection (IDS) for ICL demonstration selection (see Fig. 2 for an illustration), which can fully leverage the merits of both dimensions, *i.e.*, diversity and similarity. Intuitively, the demonstrations that are similar to the *reason* for answering a sample are strongly correlated with this sample. Therefore, we propose to incorporate zero-shot chain-of-thought reasoning (Zero-shot-CoT) into IDS to iteratively select demonstration examples that are diverse but still have a strong correlation with the test sample.

Specifically, for each test sample \hat{x}_i , IDS mainly consists of four steps:

1. We apply **Zero-shot-CoT**, *i.e.*, “Let’s think step by step.” to the test sample \hat{x}_i before selecting demonstrations to obtain a reasoning path R .
2. The **reasoning path** R is then used to select top- k (k is the number of shot) most semantically similar training examples $\{(x_1, y_1), \dots, (x_k, y_k)\}$ as few-shot demonstrations. We use Sentence-BERT (Reimers and Gurevych, 2019) to encode the reasoning path R and training examples to obtain the contextual representations and use cosine similarity to measure the similarity between representations.
3. The selected k training examples $\{(x_1, y_1), \dots, (x_k, y_k)\}$ are then prepended to the test sample \hat{x}_i for ICL. During inference, we ensure that the generated answer \hat{A} is accompanied by its corresponding reasoning

Algorithm 1 Selection process of IDS

Require: Training set $\mathcal{D}_{\text{train}}$, test set $\mathcal{D}_{\text{test}}$, LLM $_{\theta}$, number of demonstrations k , number of iterations q and answer set $\hat{A}_{\text{all}} = \emptyset$

- 1: ENCODE all samples in $\mathcal{D}_{\text{train}}$ using Sentence-BERT \triangleright **Encode training set**
- 2: **for** \hat{x}_i in $\mathcal{D}_{\text{test}}$ **do**
- 3: APPLY Zero-shot-CoT to \hat{x}_i to obtain the reasoning path R \triangleright **Zero-shot-CoT**
- 4: **for** $j = 1, \dots, q$ **do**
- 5: ENCODE R using Sentence-BERT \triangleright **Encode reasoning path**
- 6: USE R to select top- k most similar examples $S = \{(x_1, y_1), \dots, (x_k, y_k)\}$ from $\mathcal{D}_{\text{train}}$ as demonstrations \triangleright **KNN selection**
- 7: $(\hat{A}, \hat{R}) = \text{LLM}_{\theta}(S, \hat{x}_i)$ \triangleright **ICL with Zero-shot-CoT**
- 8: $R = \hat{R}, \hat{A}_{\text{all}} = \hat{A}_{\text{all}} \cup \{\hat{A}\}$ \triangleright **Update reasoning path and answer set**
- 9: **end for**
- 10: ADOPT majority voting for \hat{A}_{all} to obtain the final result \hat{A}_{final} for the test sample \hat{x}_i \triangleright **Majority voting**
- 11: **end for**

path \hat{R} through designed prompts, *e.g.*, “The response should follow the format: Sentiment: {positive or negative}\nReason: {reason}”. Note that **Zero-shot-CoT** is also applied in this step to improve the quality of generated reasoning paths. After ICL, we go back to Step 2 for *iterations* using the *new* reasoning path \hat{R} .

4. After q rounds of iterations between Step 2 and 3, we adopt **majority voting** on all \hat{A} to obtain the final result \hat{A}_{final} .

Obviously, the selected demonstration examples are strongly correlated with the original test sample, *i.e.*, achieving similarity, as they are selected by the generated reasoning paths (see Appendix A.3 for quantitative analysis). And they can be different during iterations to achieve diversity because the reasoning paths vary in different iterations. Note that there is *no* reasoning path in few-shot demonstrations (as shown in the green part in Fig. 2). The reasoning path only exists in the output of LLMs.

In addition, we illustrate the whole selection process in Alg. 1 and show the instructions and input formats of different types of tasks for ICL in Appendix A.1.

6 Experiments

In this section, we first describe the tasks and datasets, and then introduce methods compared in our work. Finally, we present the experimental results.

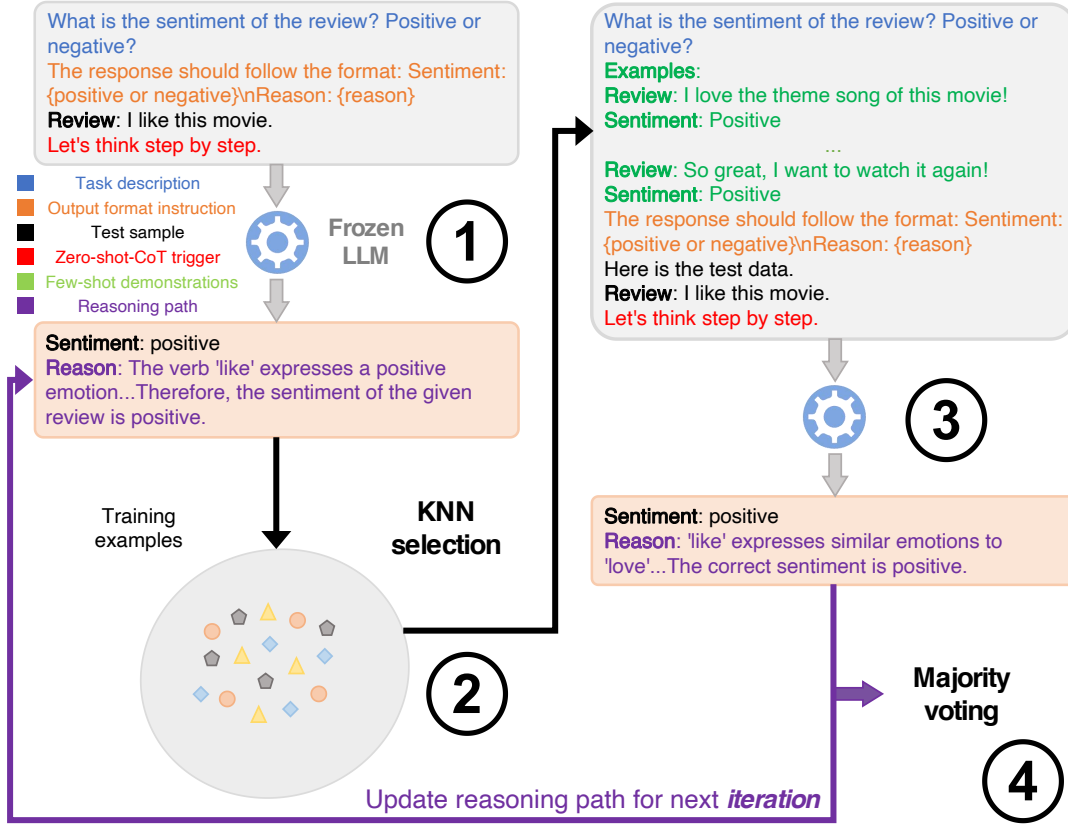


Figure 2: Illustration of our proposed Iterative Demonstration Selection (IDS). IDS first applies Zero-shot-CoT to the test sample to obtain a reasoning path, which is then used to select few-shot demonstrations from training examples through KNN. The selected demonstration examples are prepended to the test sample for ICL. To obtain the new reasoning path for extracting another set of demonstrations in the next iteration, an instruction for output format is inserted before the test sample. After several iterations, IDS uses majority voting to obtain the final result.

6.1 Experimental Setup

Tasks and Datasets We mainly investigate 6 different datasets covering 4 representative task categories: reasoning (commonsense reasoning (CommsenseQA (Talmor et al., 2019)), mathematical reasoning (GSM8K (Cobbe et al., 2021)) and logical reasoning (LogiQA (Liu et al., 2020))), question answering (BoolQ (Clark et al., 2019)), topic classification (AGNews (Zhang et al., 2015)) and sentiment analysis (SST2 (Socher et al., 2013)). For each dataset, we randomly sample at most 10000 examples from the original training set as $\mathcal{D}_{\text{train}}$ and at most 2000 test examples as $\mathcal{D}_{\text{test}}$ for evaluating the performance of selected demonstrations. The detailed information of different datasets is shown in Appendix A.2. To reduce the randomness, we run every experiment five times with different random seeds (resulting in different training and test samples if not using the whole set) and report the average results. Without specification, we use $k = 4$ number of demonstrations following Wang et al. (2023b) and set the number of iterations q to 3.

Methods Compared We mainly use GPT-3.5 (gpt-3.5-turbo) as the LLM and compare our IDS with the following methods in the experiments for selecting ICL demonstrations:

- **Top- k -Consistency** (Liu et al., 2022) selects the $top-k$ semantically similar examples from the training set $\mathcal{D}_{\text{train}}$ as demonstrations for each test sample and applies *self-consistency* (Wang et al., 2022b) with q decoding paths (temperature 0.7) to match the number of iterations. Following Zhang et al. (2023a), all samples are encoded by Sentence-BERT (Reimers and Gurevych, 2019) to obtain contextual representations for calculating the cosine similarity.
- **Random-Voting** randomly selects k examples from $\mathcal{D}_{\text{train}}$ as few-shot demonstrations for every test sample and runs experiments q times before majority voting.
- **Cluster-Voting** partitions $\mathcal{D}_{\text{train}}$ into k clusters and selects a representative example from each

Method	BoolQ	CommonsenseQA	GSM8K	LogiQA	AGNews	SST2	Average
G-fair-Prompting	84.8 \pm 0.7	75.5 \pm 0.3	76.9 \pm 0.6	43.8 \pm 0.4	88.9 \pm 1.0	94.6 \pm 0.3	77.4 \pm 0.2
Skill-KNN	85.9 \pm 0.5	75.2 \pm 0.2	76.5 \pm 0.3	44.6 \pm 0.2	88.7 \pm 0.9	94.9 \pm 0.2	77.6 \pm 0.1
Top- k -Consistency	87.1 \pm 0.2	74.5 \pm 0.2	76.1 \pm 0.5	45.7 \pm 0.4	89.3 \pm 0.8	95.2 \pm 0.4	78.0 \pm 0.1
Random-Voting	87.3 \pm 0.6	77.0 \pm 0.2	75.6 \pm 0.4	45.1 \pm 0.3	87.0 \pm 1.6	95.6 \pm 0.1	77.9 \pm 0.2
Cluster-Voting	86.4 \pm 0.7	76.5 \pm 0.3	76.8 \pm 0.3	44.1 \pm 0.3	86.8 \pm 1.2	95.2 \pm 0.4	77.6 \pm 0.2
IDS	87.8\pm0.8	78.1\pm0.1	78.5\pm0.4	46.9\pm0.2	89.8\pm0.8	95.8\pm0.2	79.5\pm0.1

Table 3: Accuracy (%) of different methods on 6 datasets. **Bold** indicates the best result. IDS is consistently better than all previous baselines.

	2	4	6	8
Top- k -Consistency	77.9	78.2	78.4	78.4
IDS	79.1	79.6	79.4	79.3

Table 4: Accuracy (%) of Top- k -Consistency and IDS with different numbers of demonstrations k .

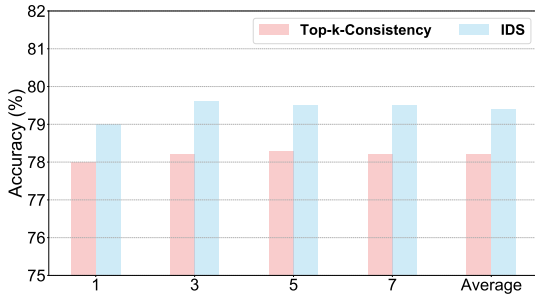


Figure 3: Accuracy (%) of Top- k -Consistency and IDS with different numbers of reasoning paths or iterations.

cluster to form demonstrations. Following Zhang et al. (2023a), we choose the sample closest to the centroid in each cluster as the representative example. Same as Random-Voting, after running experiments q times, Cluster-Voting adopts majority voting to obtain the final result.

Besides, we also compare IDS with two latest ICL demonstration selection approaches: **G-fair-Prompting** (Ma et al., 2023) and **Skill-KNN** (An et al., 2023b). Note that we find that simultaneously generating answers and reasoning paths can improve the ICL performance in general even if the target task is not a reasoning task in the conventional sense, *e.g.*, sentiment analysis. Therefore, we apply the same prompt, *e.g.*, “The response should follow the format: Sentiment: {positive or negative}\nReason: {reason}”, and *Zero-shot-CoT* to baseline methods.

6.2 Main Results

Table 3 shows the average performance scores of different methods on all investigated datasets.

From the results, we can observe that

- Our proposed IDS consistently outperforms previous baselines on all datasets with a negligible increase in API request cost (Zero-shot-CoT in the first step), which demonstrates that our method can indeed effectively and efficiently select better ICL demonstration examples. On average, IDS yields about 1.5% performance boost compared with the best baseline Top- k -Consistency as it can fully leverage the merits of both selection dimensions (diversity and similarity). While the performance gain on a few simple benchmarks looks somewhat small (because the baseline results are already pretty high, *e.g.*, the baseline performance of SST2 is above 95%), IDS performs much better than baselines on more complex tasks. For example, IDS can bring an average relative improvement of about 3% on GSM8K (mathematical reasoning) and LogiQA (logical reasoning) compared with Top- k -Consistency.

- Cluster-Voting underperforms Random-Voting and Top- k -Consistency on most datasets, which is inconsistent with the conclusion in AutoCoT (Zhang et al., 2023a). As shown in Zhang et al. (2023a), selecting a representative sample from each cluster and generating the corresponding reasoning chain using Zero-shot-CoT to construct chain-of-thought demonstrations can achieve better performance than selection by similarity or random selection. We speculate that this is because there is no rationale in ICL demonstration examples except GSM8K, which eliminates the advantage of cluster-based methods in mitigating misleading caused by rationale errors. In addition, Cluster-Voting selects demonstrations at the dataset level, *i.e.*, all test samples use the same demonstration examples, which is not as flexible as other instance-level methods.

6.3 Analysis

Different Numbers of Demonstrations While we use $k = 4$ demonstration examples for all ex-

	gpt-3.5-turbo	gpt-4
Top- k -Consistency	78.2	80.8
IDS	79.6	82.1

Table 5: Accuracy (%) of Top- k -Consistency and IDS with different LLMs (gpt-3.5-turbo and gpt-4). For gpt-4, we randomly sample 200 test examples per dataset due to the high cost.

	BoolQ	GSM8K
Top- k -Consistency	84.2	49.6
IDS	85.4	51.4

Table 6: Accuracy (%) of different methods with Llama-2-70b-chat.

periments, we also evaluate the effectiveness of IDS with different k . We randomly choose one seed for experiments and report the average results of the 6 datasets in Table 4. We can see that IDS consistently outperforms the best baseline Top- k -Consistency with different numbers of demonstrations. In addition, more demonstrations do **not** guarantee better ICL performance, which is consistent with the observation in Wang et al. (2023b).

Different Numbers of Iterations Our experiments and analysis so far use $q = 3$ iterations. To verify whether the performance gain of IDS is consistent across different numbers of iterations, we conduct controlled experiments with $q = \{1, 5, 7\}$. The average results of the 6 datasets with a randomly selected seed are reported in Fig. 3. IDS consistently outperforms Top- k -Consistency with different q . Interestingly, the performance of ICL does not always improve with the number of iterations, which might be because increased iterations can also lead to unnecessary noise.

Robustness to Model Types To demonstrate the robustness of IDS to model types, we conduct controlled experiments with GPT-4 (gpt-4). Specifically, we randomly select one seed and sample 200 test examples per dataset for experiments due to the expensive cost. From the average results reported in Table 5, we can observe that IDS still achieves better performance than Top- k -Consistency when using GPT-4 as the LLM, showing its robustness to different LLMs.

Generalization to Open-source LLMs To better verify the generalization ability of IDS, we use

Iterative Demonstration Selection Question: The homeowner frowned at the price of gas, what did he have to do later? Answer Choices: (A) own home (B) mail property tax payments (C) board windows (D) cut grass (E) receive mail Iteration 1: Answer: B\nReason: ... Iteration 2: Answer: D\nReason: ... Iteration 3: Answer: D\nReason: ... Label: D	Top-k-Consistency Question: The homeowner frowned at the price of gas, what did he have to do later? Answer Choices: (A) own home (B) mail property tax payments (C) board windows (D) cut grass (E) receive mail Response: Answer: B\nReason: ...; Answer: B\nReason: ...; Answer: B\nReason: ... Label: D
Iterative Demonstration Selection Input: Texas entrepreneur wants to kick computer gaming up to the next level by offering players a chance at some real-live killing via mouse and modem. Iteration 1 Examples: Input: Six days a week, teens crowd the Blue Screen Gaming cybercafe to hunt each other down with assault rifles inside virtual computer worlds... Topic: Technology Response: Topic: Technology\nReason: ... Iteration 2: ... Response: Topic: Technology ... Iteration 3: ... Response: Topic: Technology ... Label: Technology	Random-Voting Input: Texas entrepreneur wants to kick computer gaming up to the next level by offering players a chance at some real-live killing via mouse and modem. Iteration 1 Examples: Input: The Boston Celtics added a healthy Tom Gugliotta and deleted injured Delonte West. Tom, 34, was activated Wednesday from the injured list after missing seven games ... Topic: Sports Response: Topic: Sports\nReason: ... Iteration 2: ... Response: Topic: Business ... Iteration 3: ... Response: Topic: Sports ... Label: Technology

Figure 4: Several case studies of model responses. We color correct outputs in green, and wrong outputs in red.

vLLM (Kwon et al., 2023) to serve a Llama-2-70b-chat model (Touvron et al., 2023) for experiments and compare IDS with the best baseline Top- k -Consistency on two datasets: BoolQ and GSM8K. We randomly sample 500 test examples for experiments and report the results in Table 6, which demonstrates that IDS can successfully generalize to open-source LLMs.

Case Study To further understand the advantage of IDS, we show several cases in Fig. 4. As shown in the upper part of the figure, IDS can iteratively select more diverse demonstration examples than Top- k -Consistency which may be able to correct errors from previous iterations. Compared with Random-Voting, IDS can find examples that share more similar input-output patterns with the test sample to induce the LLM to generate correct answers (the lower part of the figure).

In addition, we show the robustness of IDS to different embedding models and Zero-shot-CoT triggers, and the analysis of average similarity scores in Appendix A.4 ~ A.6, respectively.

7 Conclusion

In this work, we have introduced Iterative Demonstration Selection (IDS) that can iteratively select examples that are diverse but still strongly correlate with the test sample as demonstrations by leveraging Zero-shot-CoT to improve the performance of in-context learning (ICL). Extensive experimental results and analysis show that IDS can consistently outperform previous ICL demonstration selection baselines.

Limitations

This work has several limitations. First, due to the inference cost of ChatGPT¹, we do not conduct experiments on the entire test set. Besides, we include 6 datasets covering 4 different task types in this work. A further improvement could be to explore more diverse types of tasks.

References

- Shengnan An, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Jian-Guang Lou, and Dongmei Zhang. 2023a. [How do in-context examples affect compositional generalization?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11027–11052, Toronto, Canada. Association for Computational Linguistics.
- Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Weizhu Chen, and Jian-Guang Lou. 2023b. [Skill-based few-shot selection for in-context learning](#). *arXiv preprint arXiv:2305.14210*.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *arXiv preprint arXiv:2005.14165*.
- Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Srinu Iyer, Veselin Stoyanov, and Zornitsa Kozareva. 2022. [Improving in-context few-shot learning via self-supervised training](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3558–3573, Seattle, United States. Association for Computational Linguistics.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [BoolQ: Exploring the surprising difficulty of natural yes/no questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2023. [Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers](#). In *ICLR 2023*

Workshop on Mathematical and Empirical Understanding of Foundation Models.

- Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. 2023. [Active prompting with chain-of-thought for large language models](#). *arXiv preprint arXiv:2302.12246*.
- Bosheng Ding, Chengwei Qin, Linlin Liu, Lidong Bing, Shafiq Joty, and Boyang Li. 2022. [Is gpt-3 a good data annotator?](#) *arXiv preprint arXiv:2212.10450*.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. [A survey for in-context learning](#). *arXiv preprint arXiv:2301.00234*.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. [Model-agnostic meta-learning for fast adaptation of deep networks](#). In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR.
- Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. 2023. [Specializing smaller language models towards multi-step reasoning](#). *arXiv preprint arXiv:2301.12726*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). 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 3816–3830, Online. Association for Computational Linguistics.
- Tianyu Gao, Xu Han, Ruobing Xie, Zhiyuan Liu, Fen Lin, Leyu Lin, and Maosong Sun. 2020. [Neural snowball for few-shot relation learning](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7772–7779. AAAI Press.
- Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022. [Large language models are reasoning teachers](#). *arXiv preprint arXiv:2212.10071*.
- Zikun Hu, Xiang Li, Cunchao Tu, Zhiyuan Liu, and Maosong Sun. 2018. [Few-shot charge prediction with discriminative legal attributes](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 487–498, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. 2023. [Language is not all you need: Aligning perception with language models](#). *arXiv preprint arXiv:2302.14045*.

¹experiments done between 06/2023 and 10/2023

648	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners . In <i>Thirty-sixth Conference on Neural Information Processing Systems (NeurIPS 2022)</i> .	705
649		706
650		707
651		708
652		709
653	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention . In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	710
654		711
655		712
656		713
657		714
658		715
659		716
660	Itay Levy, Ben Bogin, and Jonathan Berant. 2022. Diverse demonstrations improve in-context compositional generalization . <i>arXiv preprint arXiv:2212.06800</i> .	717
661		718
662		719
663		720
664	Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023a. Unified demonstration retriever for in-context learning . <i>arXiv preprint arXiv:2305.04320</i> .	721
665		
666		
667		
668	Xiaonan Li and Xipeng Qiu. 2023. Finding supporting examples for in-context learning . <i>arXiv preprint arXiv:2302.13539</i> .	722
669		723
670		724
671	Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2022. On the advance of making language models better reasoners . <i>arXiv preprint arXiv:2206.02336</i> .	725
672		726
673		727
674		728
675	Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. 2023b. Transformers as algorithms: Generalization and stability in in-context learning .	729
676		730
677		731
678		732
679	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for GPT-3? In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.	733
680		734
681		735
682		736
683		
684		
685		
686		
687	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning . <i>arXiv preprint arXiv:2007.08124</i> .	737
688		738
689		739
690		
691		
692	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 <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.	740
693		741
694		742
695		743
696		744
697		
698		
699		
700	Man Luo, Xin Xu, Zhuyun Dai, Panupong Pasupat, Mehran Kazemi, Chitta Baral, Vaiva Imbrasaitė, and Vincent Y Zhao. 2023. Dr. icl: Demonstration-retrieved in-context learning . <i>arXiv preprint arXiv:2305.14128</i> .	745
701		746
702		747
703		748
704		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761

- Sachin Ravi and Hugo Larochelle. 2017. [Optimization as a model for few-shot learning](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). 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 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- 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.
- Alexander Scarlatos and Andrew Lan. 2023. [Ret-icl: Sequential retrieval of in-context examples with reinforcement learning](#). *arXiv preprint arXiv:2305.14502*.
- Seongjin Shin, Sang-Woo Lee, Hwiyeon Ahn, Sungdong Kim, HyoungSeok Kim, Boseop Kim, Kyunghyun Cho, Gichang Lee, Woomyoung Park, Jung-Woo Ha, and Nako Sung. 2022. [On the effect of pretraining corpora on in-context learning by a large-scale language model](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5168–5186, Seattle, United States. Association for Computational Linguistics.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. [Reflexion: Language agents with verbal reinforcement learning](#).
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, A. Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Conference on Empirical Methods in Natural Language Processing*.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. [Black-box tuning for language-model-as-a-service](#). In *International Conference on Machine Learning*, pages 20841–20855. PMLR.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *arXiv preprint arXiv:2307.09288*.
- Eleni Triantafillou, Richard Zemel, and Raquel Urtasun. 2017. [Few-shot learning through an information retrieval lens](#). *arXiv preprint arXiv:1707.02610*.
- Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. 2023a. [Neural codec language models are zero-shot text to speech synthesizers](#). *arXiv preprint arXiv:2301.02111*.
- Xinyi Wang, Wanrong Zhu, and William Yang Wang. 2023b. [Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning](#). *arXiv preprint arXiv:2301.11916*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2022a. [Rationale-augmented ensembles in language models](#). *arXiv preprint arXiv:2207.00747*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2022b. [Self-consistency improves chain of thought reasoning in language models](#). *arXiv preprint arXiv:2203.11171*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022c. [Self-instruct: Aligning language model with self generated instructions](#). *arXiv preprint arXiv:2212.10560*.
- Zhendong Wang, Yifan Jiang, Yadong Lu, Yelong Shen, Pengcheng He, Weizhu Chen, Zhangyang Wang, and Mingyuan Zhou. 2023c. [In-context learning unlocked for diffusion models](#). *arXiv preprint arXiv:2305.01115*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. [Emergent abilities of large language models](#). *Transactions on Machine Learning Research*. Survey Certification.
- 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](#). In *Thirty-sixth Conference on Neural Information Processing Systems (NeurIPS 2022)*.
- Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu, Denny Zhou, Tengyu Ma, et al. 2023a. [Symbol tuning improves in-context learning in language models](#). *arXiv preprint arXiv:2305.08298*.

875	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and	927
876	Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,	Sameer Singh. 2021. Calibrate before use: Improv-	928
877	Da Huang, Denny Zhou, et al. 2023b. Larger	ing few-shot performance of language models . In	929
878	language models do in-context learning differently .	<i>Proceedings of the 38th International Conference</i>	930
879	<i>arXiv preprint arXiv:2303.03846</i> .	<i>on Machine Learning</i> , volume 139 of <i>Proceedings</i>	931
		<i>of Machine Learning Research</i> , pages 12697–12706.	932
880	Chenfei Wu, Shengming Yin, Weizhen Qi, Xi-	PMLR.	933
881	aodong Wang, Zecheng Tang, and Nan Duan.		
882	2023. Visual chatgpt: Talking, drawing and edit-	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong	934
883	ing with visual foundation models . <i>arXiv preprint</i>	Wu, Jingjing Xu, and Baobao Chang. 2023. Can we	935
884	<i>arXiv:2303.04671</i> .	edit factual knowledge by in-context learning? <i>arXiv</i>	936
		<i>preprint arXiv:2305.12740</i> .	937
885	Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Ling-		
886	peng Kong. 2022. Self-adaptive in-context learning .	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,	938
887	<i>arXiv preprint arXiv:2212.10375</i> .	Nathan Scales, Xuezhi Wang, Dale Schuurmans,	939
		Olivier Bousquet, Quoc Le, and Ed Chi. 2022.	940
888	Sang Michael Xie, Aditi Raghunathan, Percy Liang,	Least-to-most prompting enables complex reason-	941
889	and Tengyu Ma. 2022. An explanation of in-context	ing in large language models . <i>arXiv preprint</i>	942
890	learning as implicit bayesian inference . In <i>Interna-</i>	<i>arXiv:2205.10625</i> .	943
891	<i>tional Conference on Learning Representations</i> .		
892	Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu,	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	944
893	and Lingpeng Kong. 2023. Compositional ex-	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	945
894	emplars for in-context learning . <i>arXiv preprint</i>	Ba. 2023. Large language models are human-level	946
895	<i>arXiv:2302.05698</i> .	prompt engineers . In <i>The Eleventh International</i>	947
		<i>Conference on Learning Representations</i> .	948
896	Xi Ye, Srinivasan Iyer, Asli Celikyilmaz, Ves Stoy-	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and	949
897	anov, Greg Durrett, and Ramakanth Pasunuru. 2022.	Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing	950
898	Complementary explanations for effective in-context	vision-language understanding with advanced large	951
899	learning . <i>arXiv preprint arXiv:2211.13892</i> .	language models . <i>arXiv preprint arXiv:2304.10592</i> .	952
900	Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyun-	A Appendix	953
901	soo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee,		
902	and Taeuk Kim. 2022. Ground-truth labels matter: A	A.1 Instructions and Input Formats of	954
903	deeper look into input-label demonstrations . In <i>Pro-</i>	Different Tasks	955
904	<i>ceedings of the 2022 Conference on Empirical Meth-</i>		
905	<i>ods in Natural Language Processing</i> , pages 2422–	We show the instructions and input formats of	956
906	2437, Abu Dhabi, United Arab Emirates. Association	different types of tasks for in-context learning in	957
907	for Computational Linguistics.	Fig. 5 .	958
908	Eric Zelikman, Yuhuai Wu, and Noah D Goodman.		
909	2022. Star: Bootstrapping reasoning with reason-	A.2 Datasets Information	959
910	ing . <i>arXiv preprint arXiv:2203.14465</i> .		
911	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	We show the detailed information of different	960
912	Character-level convolutional networks for text classi-	datasets in Table 7 .	961
913	fication . <i>Advances in neural information processing</i>		
914	<i>systems</i> , 28:649–657.	A.3 Measure of Reasoning Path Correlation	962
915	Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Ac-		
916	tive example selection for in-context learning . <i>arXiv</i>	We report the average similarity score between test	963
917	<i>preprint arXiv:2211.04486</i> .	samples and the corresponding generated reasoning	964
		paths ($\text{score}_{\text{reason}}$), the average similarity score be-	965
918	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex	tween test samples and randomly selected training	966
919	Smola. 2023a. Automatic chain of thought prompt-	examples ($\text{score}_{\text{random}}$), and the average similarity	967
920	ing in large language models . In <i>The Eleventh In-</i>	score between test samples and the most similar	968
921	<i>ternational Conference on Learning Representations</i>	training examples ($\text{score}_{\text{similar}}$) in Table 8 . For each	969
922	<i>(ICLR 2023)</i> .	dataset, we randomly select 100 test samples and	970
923	Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao,	use Sentence-BERT for similarity calculation. We	971
924	George Karypis, and Alex Smola. 2023b. Multi-	can observe that $\text{score}_{\text{reason}}$ is slightly worse than	972
925	modal chain-of-thought reasoning in language mod-	$\text{score}_{\text{similar}}$ and much higher than $\text{score}_{\text{random}}$, indi-	973
926	els . <i>arXiv preprint arXiv:2302.00923</i> .	cating that the generated reasoning path is indeed	974
		strongly correlated with the test sample.	975

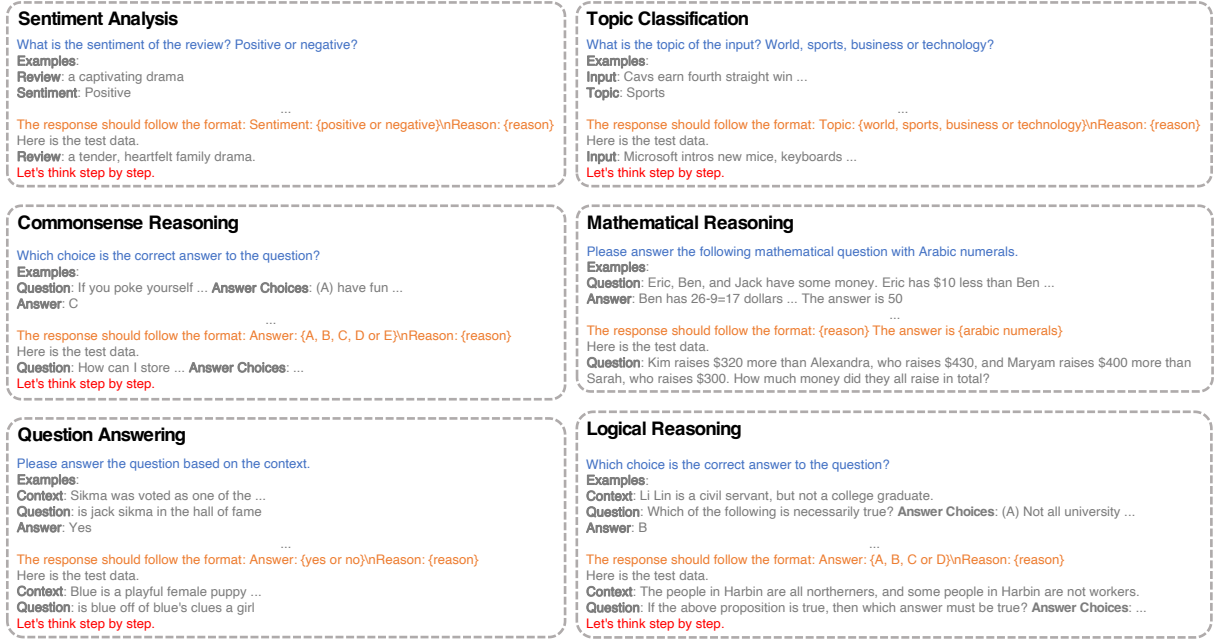


Figure 5: Instructions and input formats of four different categories of tasks (sentiment analysis, topic classification, reasoning, and question answering) for ICL. For Zero-shot-CoT in the first step of IDS, there is no demonstration example and the instruction “Here is the test data.”.

	BoolQ	CommonsenseQA	GSM8K	LogiQA	AGNews	SST2
# Training Samples	9427 (full)	9741 (full)	7473(full)	7376(full)	10000	10000
# Test Samples	2000	1221 (full)	1000	500	1000	872 (full)

Table 7: Detailed information of different datasets. # refers to ‘the number of’ and ‘full’ means the whole set. Note that different random seeds do not result in different samples if the whole set is used.

	score _{reason}	score _{random}	score _{similar}
Average Similarity Score	0.60	0.32	0.69

Table 8: Comparison between different average similarity scores.

	BoolQ	CommonsenseQA	GSM8K
Top- <i>k</i> -Consistency	86.0	75.4	75.8
IDS	87.2	78.0	77.6

Table 9: Accuracy (%) of different methods with OpenAI embedding model (text-embedding-ada-002) on three datasets.

A.4 Robustness to Embedding Models

Instead of using Sentence-BERT, we also explore adopting the OpenAI embedding model (text-embedding-ada-002) as the encoder. Specifically, we conduct experiments on 3 datasets: BoolQ, CommonsenseQA and GSM8K. For each dataset, we randomly sample 500 test examples and compare IDS with the baseline Top-*k*-Consistency. The

	Default	Trigger1	Trigger2
IDS	79.8	80.0	79.7

Table 10: Accuracy (%) of IDS with different Zero-shot-CoT triggers.

results reported in Table 9 demonstrate IDS’s robustness to different embedding models.

A.5 Robustness to Zero-shot-CoT Triggers

To verify the robustness of IDS to Zero-shot-CoT triggers, we conduct controlled experiments with two new triggers: “Let’s work this out in a step by step way to be sure we have the right answer.” (Trigger1) and “Let’s solve this problem step by step” (Trigger2). Specifically, we randomly sample 500 test examples per dataset for experiments and report the average results in Table 10, which demonstrates that IDS is indeed robust to different Zero-shot-CoT triggers.

	Top- k -Consistency	IDS	Random-Voting
Average Similarity Score	0.69	0.46	0.31

Table 11: Average similarity scores between test examples and the corresponding selected demonstrations of three methods (Top- k -Consistency, IDS and Random-Voting).

A.6 Average Similarity Scores

In Table 11, we report the average similarity scores between test samples and the corresponding demonstrations of different methods. Specifically, we randomly select 200 test examples for each dataset and use Sentence-BERT to obtain contextual representations for calculating similarity scores. We can see that the average similarity score of IDS is between that of Top- k -Consistency and Random-Voting, indicating that it can indeed strike a balance between two selection dimensions.