

# The role of diversity in in-context learning for large language models

Anonymous ACL submission

## Abstract

In-context learning (ICL) is a crucial capability of current large language models (LLMs), where the selection of examples plays a key role in performance. While most existing approaches focus on selecting the most *similar* examples to the query, the impact of *diversity* in example selection remains underexplored. We systematically investigate the role of *diversity* in in-context example selection through experiments across a range of tasks, from sentiment classification to more challenging math and code problems. Experiments on Llama-3.1, Gemma-2, and Mistral-v0.3 families of models show that diversity-aware selection methods improve performance, particularly on complex tasks like math and code, and enhance robustness to out-of-distribution queries. To support these findings, we introduce a theoretical framework that explains the benefits of incorporating diversity in in-context example selection.

## 1 Introduction

In-context learning (ICL) (Brown et al., 2020) has emerged as one of the most significant and versatile capabilities of large language models (LLMs). This paradigm allows a model to adapt to a vast array of new tasks on the fly, simply by conditioning on a prompt containing a few input-output examples, all without requiring updates to its parameters. The power and resource efficiency of ICL have made it a cornerstone of LLM applications, ranging from simple text classification (Min et al., 2022) and commonsense reasoning (Srivastava et al., 2023) to complex, multi-step tasks like mathematical problem-solving (Wei et al., 2022) and code generation (Chen et al., 2021).

The effectiveness of ICL, however, is highly sensitive to the choice of in-context examples (Lu et al., 2021; Liu et al., 2021; Chang and Jia, 2023). This makes example selection a critical area of study. To address this, prior work has explored various selection strategies: choosing examples most *similar* to the query in embedding space (Liu et al.,

2021; Yang et al., 2022; Wu et al., 2023; Qin et al., 2023), maximizing feature *coverage* (Levy et al., 2023; Ye et al., 2023; Gupta et al., 2023), selecting based on *difficulty* (Ma et al., 2025; Swayamdipta et al., 2020; Yuan et al., 2025; Cook et al., 2025), or choosing examples based on *sensitivity* (Chen et al., 2023). Other approaches train deep neural retrievers (Karpukhin et al., 2020; Rubin et al., 2022; Luo et al., 2023; Scarlatos and Lan, 2023) or leverage feedback from large language models to guide selection (Li and Qiu, 2023a; Chen et al., 2023; Wang et al., 2023). More discussion can be found in Appendix A. Among these, similarity-based methods remain the fundamental baseline due to their conceptual simplicity and consistent empirical success. However, relying solely on similarity can lead to redundancy among demonstrations and potentially omit important but less similar features (Levy et al., 2023; Gupta et al., 2023).

Within machine learning, *diversity* is also a fundamental principle for building robust and generalizable models, and its importance is widely recognized in related domains—such as fixed-prompt ICL with global demonstration sets (Li and Qiu, 2023b; Luo et al., 2024), active learning (Giouroukis et al., 2025; Shi and Shen, 2016), coreset construction (Wan et al., 2024; Zhan et al., 2025; Sener and Savarese, 2018), and instruction tuning (Wang et al., 2024). By exposing a model to a varied set of examples, we can prevent overfitting and encourage the learning of more abstract, transferable patterns. Given its foundational role, a deep understanding of diversity is crucial for unlocking the full potential of in-context learning.

Despite its importance, the explicit role of diversity in retrieval-based ICL remains underexplored. While some recent work has incorporated feature coverage as a proxy for diversity (Levy et al., 2023; Ye et al., 2023; Gupta et al., 2023), this approach is limited in scope. Coverage primarily aims to span the concrete input features of a given query, which

is a narrower goal than promoting the broader representational variety that is central to true diversity; see also examples in Appendix E that compare coverage with diversity. Other diversity-aware approaches have also been proposed, such as the S3 method from Kumari et al. (2024). However, that study was confined to simple classification and selection tasks. Consequently, the effectiveness of diversity in more complex, reasoning-based ICL applications remains an open question.

Furthermore, pursuing diversity without care can be counterproductive, as explicit diversity-aware selection risks retrieving examples that are too dissimilar from the query, potentially harming performance (An et al., 2023a). The field, therefore, lacks a systematic understanding of this trade-off. It remains unclear whether and when explicit diversity is beneficial—especially for tasks that lack clear local structure and demand more abstract reasoning. This gap in knowledge motivates the following fundamental questions:

*Should we explicitly consider diversity when selecting in-context examples? If so, under what conditions does it outperform similarity-based methods? And fundamentally, why does diversity help?*

## 1.1 Our contributions

We present, to the best of our knowledge, the first systematic investigation of the role of diversity in in-context learning. Our study spans a broad range of tasks—including sentiment classification, commonsense reasoning, math generation, reading comprehension, and SQL code generation—covering diverse types and varying levels of difficulty. We compare five demonstration selection methods: (1) random selection (Rand); (2) selecting the  $K$  most similar examples to the query (TopK) (Liu et al., 2021); (3) Select the most representative examples from a similarity-reduced subset (Div-S3) (Kumari et al., 2024); (4) selecting similar examples from a diversity-reduced subset (Div) (Su et al., 2023), which relates to DPP-based diversity (Chen et al., 2018); and (5) a sequential method that balances similarity to the query and diversity among selected examples (TopK-Div). We are the first to systematically evaluate methods for (4) and (5) in the ICL setting. Their approaches are particularly compelling because they offer explicit control over the diversity level, enabling a tunable trade-off between selecting highly relevant examples and avoiding redundancy—a key factor

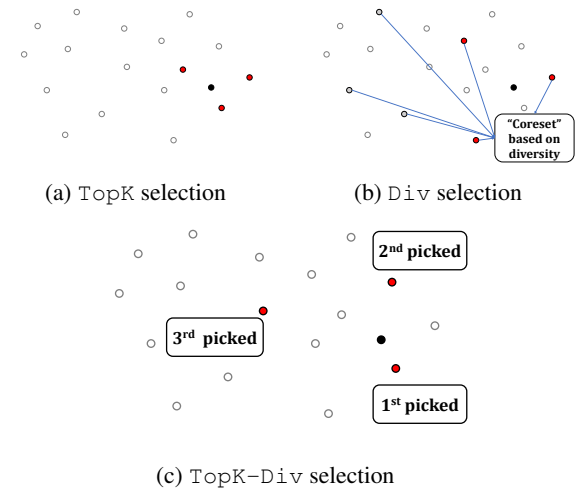


Figure 1: An illustrative example for TopK, Div, and TopK-Div methods. Point filled in black denotes the query. (a) TopK: Select the most similar demonstrations (red points) in the embedding space. (b) Div: First select a “coreset” based on some diversity metric, which is fixed for all queries (red/gray points). Then select the most similar demonstrations (red points) among this “coreset”. (c) TopK-Div: Select the demonstrations sequentially based on the linear combination of similarity to the query and the diversity with the selected examples. The first example is the one most similar to the query. When picking the second, it is balanced between the similarity to the query and the diversity from the first example.

in optimizing LLM performance within limited context windows. See Figure 1 for illustration and Section 2 for formal definitions.

Through experiments on frontier open-source models (Llama-3.1 (Dubey et al., 2024), Gemma-2 (Team et al., 2024), and Mistral-v0.3 (Jiang et al., 2023)), we reach the following findings. **Finding 1:** Diversity-aware demonstration selection methods achieve better performance on more “challenging” tasks like reading comprehension, math, and code. As task difficulty increases, diversity-aware methods yield greater relative benefits, narrowing the gap with the TopK method or even surpassing it.

While changing the tasks and even the language model to use will change the ranking of the demonstration selection methods we test, in general we find that diversity-aware methods, namely Div, TopK-Div and Div-S3, perform better on more challenging tasks like reading comprehension, math, and simple code generation. On the other hand, on simple tasks like sentiment classification and multiple-choice, TopK performs the best (Table 1). TopK-Div achieves, on average, more than

160 a 1% improvement over TopK on difficult tasks, 211  
 161 whereas on simpler tasks TopK holds a marginal 212  
 162 0.14% average advantage over TopK-Div. 213  
 163 Quantitative analysis of performance improvements under 214  
 164 varying levels of added diversity for TopK-Div 215  
 165 and Div demonstrates that more challenging tasks 216  
 166 benefit more from increased diversity, further vali- 217  
 167 dating this finding (Figure 4, Table 7). 218

168 **Finding 2:** *Diversity-aware methods work bet- 219*  
 169 *ter for out-of-distribution queries.* When the 220  
 170 query and demonstrations come from different 221  
 171 distributions, diversity-aware methods are more 222  
 172 likely to perform well. For example, on senti- 223  
 173 ment classification, when both demonstration and 224  
 174 query come from the SST-2 dataset, which con- 225  
 175 sists of movie reviews, TopK performs the best, 226  
 176 and there is a gap with all other methods (Table 2). 227  
 177 The average performance gap between TopK and 228  
 178 TopK-Div is 1.26% across the two models. How- 229  
 179 ever, when changing the demonstrations from SST- 230  
 180 2 to IMDB (which also consists of movie reviews), 231  
 181 TopK-Div outperforms TopK by 0.6% on average 232  
 182 (Table 2). A similar observation holds for vari- 233  
 183 ous splits of Geoquery dataset (Figure 2). 234

184 **Finding 3:** In the same task, diversity-aware 235  
 185 methods likely perform better on “harder” exam- 236  
 186 ples, e.g. reading comprehension with longer con- 237  
 187 text, or SQL code generation with more structures 238  
 188 (Table 3). On the easier samples from GeoQuery 239  
 189 and SQuAD, TopK-Div achieves an average accu- 240  
 190 racy improvement of 2.12% over TopK. On the 241  
 191 more challenging samples, the average improve- 242  
 192 ment of TopK-Div over TopK increases to 6.47%. 243  
 193 Our analysis also reveals that diversity exhibits a 244  
 194 “beyond-coverage” phenomenon, both at the task 245  
 195 level and the example level (see discussion in Sec- 246  
 196 tion 3.1 and Appendix E). 247

197 We discuss these findings in detail in Section 3. 248  
 198 To ensure robustness, we perform extensive abla- 249  
 199 tions across model scales (1B–70B) and in-context 250  
 200 demonstration counts. Beyond empirical trends, 251  
 201 our study provides actionable guidance for tuning 252  
 202 demonstration diversity across tasks such as reason- 253  
 203 ing and generation. Overall, our findings deepen 254  
 204 the understanding of how diversity influences in- 255  
 205 context learning, and inform principled strategies 256  
 206 for demonstration selection in real-world applica- 257  
 207 tions. 258

## 208 2 Background and notations 259

209 In this section, we introduce the in-context learning 260  
 210 (ICL) paradigm, relevant demonstration selection 261

211 methods, and associated notations. 212

213 **In-context learning (ICL).** A task  $\mathcal{T} =$  214  
 215  $(\mathcal{X}, \mathcal{Y}, P(y|x))$  defines a probabilistic mapping 216  
 217 from an input  $x \in \mathcal{X}$  to an output  $y \in \mathcal{Y}$ . For exam- 218  
 219 ple, the task can be sentiment classification where 220  
 221 the input space contains reviews of products and the 222  
 223 output space contains the customer’s correspond- 224  
 225 ing sentiment (positive or negative). We are pro- 226  
 227 vided with a demonstration set  $D = \{(x_i, y_i)\}_{i=1}^n$ , 228  
 229 where inputs  $x_i$  are drawn from a demonstration 230  
 231 distribution  $\mathcal{D}_{\mathcal{X}}$  and  $y_i \sim P_{\mathcal{T}}(y|x_i)$ . Queries  $x_q$  232  
 233 are drawn from a query distribution  $\mathcal{Q}_{\mathcal{X}}$ , which 234  
 235 may differ from  $\mathcal{D}_{\mathcal{X}}$  (representing shifts in domain 236  
 237 or complexity). For math tasks, the demonstra- 238  
 239 tion set may contain many elementary-level prob- 240  
 241 lems, while the query may require solving more ad- 241  
 242 vanced, middle-school-level problems. Now given 242  
 243 a query input  $x_q \sim \mathcal{Q}_{\mathcal{X}}$ , the in-context learning 243  
 244 paradigm refers to the following capability of a 244  
 245 large language model. 245

246 **Definition 2.1 (In-Context learning (ICL)).** 247  
 248 Given an LLM, a prompting strategy  $\text{Prompt}$ , 248  
 249 a demonstration set  $D = \{(x_i, y_i)\}_{i=1}^n$ , and a 249  
 250 query  $x_q$ , ICL involves selecting a small subset 250  
 251  $S = \{(x_{j_i}, y_{j_i})\}_{i=1}^K$  with shots  $K$  from the demon- 251  
 252 strations  $D$ . The LLM then predicts the output  $y_q$  252  
 253 for  $x_q$  as:  $P_{\mathcal{T}}(y|x_q) \approx \text{LLM}(\text{Prompt}(S, x_q))$ . 253

254 **Demonstration selection for ICL.** Choosing a 254  
 255 small subset  $S$  (see Theorem 2.1) is vital due to 255  
 256 LLM context limits, efficiency needs, and the ob- 256  
 257 servation that excessive demonstrations can impair 257  
 258 performance. Prior work has shown that ICL per- 258  
 259 formance is highly sensitive to this selection (Liu 259  
 260 et al., 2021), and thus sparks the study for *demon- 260*  
 261 *stration selection*. While numerous selection strate- 261  
 262 gies are proposed, the most notable and effective 262  
 263 methods are the ones that select the demonstrations 263  
 264 most similar to the query in the embedding space. 264  
 265 Efforts are also made to retrieve the demonstrations 265  
 266 using another model (can be another LLM), as well 266  
 267 as considering diversity/coverage. However, there 267  
 268 is no consensus on which method to use in a spec- 268  
 269 ific setting, and there is nearly no understanding of 269  
 270 these methods (further discussed in Appendix A). 270

271 To analyze these methods and the role of di- 271  
 272 versity, we focus on five representative selection 272  
 273 strategies: 273

274 **Method 1: Rand.** For a query  $x_q$ , this method 274  
 275 uniformly and randomly selects  $K$  demonstrations 275  
 276 from the set  $D$ . Note that Rand can also be viewed 276  
 277 as a method that is aware of diversity, but it has 277

nothing to do with the coverage.

**Method 2: TopK.** This method selects  $K$  demonstrations from  $D$  that exhibit the highest cosine similarity to the query  $x_q$  within an embedding space mapped by  $E : \mathcal{X} \rightarrow \mathcal{E}$ . It maximizes

$$\text{Sim}(E(x_i), E(x_q)) := \frac{\langle E(x_i), E(x_q) \rangle}{\|E(x_i)\| \cdot \|E(x_q)\|} \quad (1)$$

**Method 3: Div-S3.** This method, proposed by Kumari et al. (2024) for in-context demonstration selection, combines a similarity-based pruning step with a greedy submodular optimization to select examples that are both relevant and diverse. The approach aims to ensure representative coverage while maintaining closeness to the query. Although submodular diversity techniques have been well-studied in classical data selection (Lin and Bilmes, 2011; Prasad et al., 2014), their application in ICL has not been systematically explored.

**Method 4: Div.** This approach first constructs a diverse “coreset”  $D_r \subset D$  of size  $m$  (where  $K \leq m \leq n$ ). Starting with one randomly chosen demonstration,  $D_r$  is built greedily by adding  $(x, y) \in D \setminus D_r$  that maximizes

$$\text{Div}(E(x), D_r) := 1 - \frac{\sum_{(x_j, y_j) \in D_r} \text{Sim}(E(x), E(x_j))}{|D_r|} \quad (2)$$

we stop after  $D_r$  contains  $m$  examples. This is the procedure to get a diverse set of demonstrations for a task (Su et al., 2023). Subsequently, TopK selection is applied to  $D_r$  to choose  $K$  demonstrations for the query  $x_q$ . The coreset size  $m$  controls the trade-off between diversity and similarity: setting  $m = K$  emphasizes diversity by forcing selection from a small pool, while increasing  $m$  shifts the method closer to TopK by enlarging the candidate set based on similarity.

**Method 5: TopK-Div.** This method serves as a combination of TopK and Div, which includes some awareness of the diversity through similarity-based selection. It is also a greedy-like procedure when selecting the demonstration set  $S$ . Suppose that  $S$  does not reach size  $K$ , then we select the demonstration  $(x, y) \in D \setminus K$  that maximize the following metric:

$$\alpha \cdot \text{Sim}(E(x), e_q) + (1 - \alpha) \cdot \text{Div}(E(x), S) \quad (3)$$

The hyperparameter  $\alpha$  governs the balance between diversity and similarity: setting  $\alpha = 0$  emphasizes diversity among selected examples, while  $\alpha = 1$  recovers the TopK method that prioritizes similarity to the query. We stop when  $S$  has size

$K$ . For the first demonstration (when  $S$  is empty),  $\text{Diversity}(E(x), S)$  is defined as 0, thus prioritizing similarity.

The use of TopK-Div and Div methods for demonstration selection is, to our knowledge, new in the ICL setting. Their flexibility in adjusting the diversity level offers practical value, as it enables task-specific tuning to improve performance; see Section 3.1 for details.

### 3 Experiments and findings

This section empirically tests whether diversity-aware retrieval (Div, TopK-Div) yields more reliable in-context learning than similarity-only baselines (TopK).

**Tasks and datasets.** We consider 5 tasks: sentiment classification (classification task), commonsense reasoning (multiple-choice), text to SQL generation (generation), math (generation), and reading comprehension (generation). For sentiment classification, we test on SST-2 (Scarlatos and Lan, 2023), IMDB (Maas et al., 2011) and Amazon (polarity) (McAuley and Leskovec, 2013). For commonsense reasoning, we use ARC-Easy (Clark et al., 2018) and CommonsenseQA (CsQA) (Talmor et al., 2019). For text to SQL generation, we use GeoQuery (Zelle and Mooney, 1996; Tang and Mooney, 2001). For math problems, we test on GSM8K (Cobbe et al., 2021) and GSM-Plus-Mini (Li et al., 2024) datasets. For reading comprehension, we use SQuAD (Rajpurkar et al., 2016) and SCIQ (Welbl et al., 2017) datasets. We subsample some datasets to reduce the computation resources needed.

**Models.** Our main experiments are conducted on Llama 3.1 and Llama 3.2 (Dubey et al., 2024), Gemma 2 (Team et al., 2024), and Mistral v0.3 (Jiang et al., 2023) families of models. For math problems (GSM8K and GSM-Plus-Mini), we use the instruction-tuned LLMs, while for all other datasets, we use the base model. For the main experiments, we use Sentence-BERT (Reimers, 2019) to compute all the embeddings for TopK, Div, and TopK-Div. Experiments are conducted on 2 A100 GPUs.

**Hyperparameters.** For Div, we choose to first reduce the demonstration set  $D$  to a “coreset”  $D_r$  with size 100. This choice balances full similarity selection (TopK) and methods focusing mainly on diversity. For TopK-Div, we choose  $\alpha = 1/2$  to balance between similarity and diversity. For Div-S3, we choose  $|D_r| = 100$ . For the classifi-

cation task, we predict positive if the logit for token `great` is larger than that for token `terrible` for the next token prediction given a prompt. For multiple-choice tasks, we choose the option with the lowest average CE loss given a prompt. For generation tasks (text to SQL, math, reading comprehension), we use greedy decoding. More experimental details, including the prompt for each task, can be found in Appendix B.

### 3.1 Main findings

**Finding 1: Diversity-aware methods perform better on more “challenging” tasks.** Table 1 summarizes our main results in the in-distribution (ID) setting, where the demonstration distribution  $\mathcal{D}_X$  matches the query distribution  $\mathcal{Q}_X$ . For simpler tasks like sentiment classification, TopK consistently performs best, significantly outperforming diversity-emphasizing methods like Rand and Div (e.g., TopK outperforms these methods by at least 1% on average in SST-2), while TopK-Div (balancing similarity and diversity) typically ranks between these extremes. In commonsense reasoning (multiple-choice), introducing some diversity via TopK-Div improves performance over pure TopK, as observed in Commonsense QA, although the gains on ARC-Easy remain modest.

For more complex tasks—including reading comprehension, text-to-SQL generation, multi-step mathematical reasoning, and GeoQuery—introducing diversity consistently improves performance over the similarity-based TopK baseline. In GeoQuery specifically, diversity (TopK-Div) yields at least a 7% absolute accuracy gain, likely due to enhanced feature coverage (Levy et al., 2023; Ye et al., 2023).<sup>1</sup> However, excessive diversity (Rand and Div) becomes detrimental, as overly dissimilar examples fail to illustrate coherent solution patterns.

For math and reading comprehension tasks, methods that emphasize diversity—such as Div and even Rand—outperform TopK. Interestingly, the effectiveness of random selection cannot be explained solely by coverage, as random demonstrations do not systematically capture similar problem structures. Instead, we observe that for tasks where the model already exhibits strong zero-shot abilities (e.g., Math and Reading), incorporating

<sup>1</sup>The correspondence between inputs and outputs is deterministic, and the model need to learn this mapping from the provided examples. Better coverage of the query inputs implies that the model acquires a larger portion of the mappings required for the query inputs.

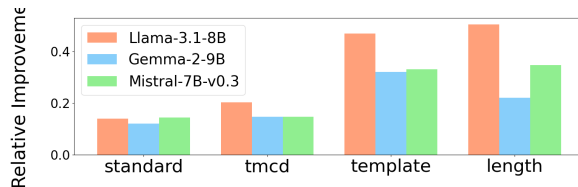


Figure 2: **(Comparison of different methods on GeoQuery OOD setting)** We report the relative improvement of TopK-Div over TopK when demonstrations and queries come from different GeoQuery dataset splitting ways. “standard” split denotes ID the setting. The relative improvement enlarges in the OOD setting.

diverse demonstrations encourages the model to rely more on its general reasoning skills rather than memorizing surface-level patterns. To support this interpretation, we present 0-shot and 1-shot performance results in Appendix C.1, which highlight the model’s underlying capabilities.

Regarding Div-S3, we observe that it underperforms TopK on relatively simple classification tasks, but shows a performance advantage on more complex tasks such as Math and Reading. These trends are consistent with other diversity-based methods and complement the analysis in (Kumari et al., 2024), extending their findings to a broader range of task types and difficulty levels.

To further verify the impact of diversity on different tasks, we also conducted experiments on both Div and TopK-Div by varying their degree of diversity (Figure 4 and Table 7). We observe a consistent pattern across both Div and TopK: for simpler tasks, introducing less diversity (i.e., employing larger subset size  $m$  for Div or higher  $\alpha$  for TopK-Div) leads to better performance, whereas for more complex tasks, incorporating greater diversity yields superior performance. Due to space limit, we defer a more detailed discussion to Appendices C.2 and C.3.

**Finding 2: Diversity helps out-of-distribution generalization.** Table 2 presents results on sentiment classification, commonsense reasoning, and reading comprehension, while Figure 2 shows text-to-SQL generation performance in the out-of-distribution (OOD) setting, where the demonstration distribution  $\mathcal{D}_X$  and query distribution  $\mathcal{Q}_X$  differ.

Overall, diversity improves OOD in-context learning. In sentiment classification, TopK performs best when both demonstrations and queries come from SST-2. However, when demonstrations shift to IMDB (another movie review dataset), TopK and TopK-Div perform similarly. When

Table 1: **(Comparison of different in-context example selection methods)** We compare diversity-aware methods `Div` and `TopK-Div` with randomly chosen (`Rand`) and similarity-based method `TopK` on a variety of tasks using different models with different number of in-context examples  $K$ . For `TopK` and `TopK-Div`, we test ten different permutations of the demonstration due to the determined choice by these methods; For `Rand` and `Div`, we test ten different random seeds. We use the corresponding instruct-tuned model for math tasks (`GSM8K` and `GSM-Plus-Mini`) and base model for all other tasks. For `TopK` and `TopK-Div` methods - both being deterministic approaches - we computed outcomes across ten distinct example permutations. For `Rand` and `Div` methods, we report the averaged results across ten random seeds. There is a huge improvement when the shot number increases from 0 to 4 / 8, which demonstrates the effectiveness of our example selection. Due to the absence of prior knowledge for `Geoquery` in the zero-shot ( $k = 0$ ) setting, we omit its  $k = 0$  results. The bold entries indicate optimal performances. The std is no more than 1% in most cases; see Appendix D for details.

| Model           | $K$        | Method   | Classification |              | Multiple-choice |              | Math         | Code         | Reading      |               |              |
|-----------------|------------|----------|----------------|--------------|-----------------|--------------|--------------|--------------|--------------|---------------|--------------|
|                 |            |          | SST-2          | Amazon       | ARC-Easy        | CsQA         |              |              | GSM8K        | GSM-Plus-Mini | GeoQuery     |
| Llama-3.1-8B    | 0          | -        | 87.50          | 95.40        | 82.43           | 62.80        | 53.45        | 65.12        | —            | 42.30         | 36.40        |
|                 |            | Rand     | 91.31          | 96.38        | 84.72           | 71.15        | 82.24        | 66.90        | 12.50        | 75.95         | 74.00        |
|                 |            | TopK     | <b>94.13</b>   | 96.24        | <b>86.10</b>    | 72.54        | 81.99        | 65.30        | 62.79        | 73.51         | 72.70        |
|                 | 4          | Div-S3   | 92.89          | <b>96.81</b> | 85.81           | 72.28        | <b>82.52</b> | 66.97        | 34.54        | <b>77.95</b>  | <b>74.61</b> |
|                 |            | Div      | 91.50          | 96.18        | 85.06           | 71.17        | 82.14        | <b>66.92</b> | 33.79        | 75.66         | 74.47        |
|                 |            | TopK-Div | 92.75          | 96.43        | 85.83           | <b>72.57</b> | 81.74        | 66.12        | <b>71.14</b> | 73.28         | 73.87        |
|                 | 8          | Rand     | 92.27          | 96.63        | 84.38           | 72.23        | 82.81        | 66.72        | 23.11        | 77.13         | 74.65        |
|                 |            | TopK     | 93.64          | 96.12        | <b>85.91</b>    | <b>73.91</b> | 82.26        | 65.99        | 72.04        | 75.52         | 74.72        |
|                 |            | Div-S3   | <b>93.65</b>   | <b>96.74</b> | 85.50           | 73.04        | <b>83.00</b> | <b>66.82</b> | 43.61        | <b>79.41</b>  | 75.15        |
|                 |            | Div      | 92.95          | 96.25        | 84.97           | 72.77        | 82.98        | 66.56        | 38.61        | 77.71         | <b>75.17</b> |
|                 |            | TopK-Div | 93.33          | 96.57        | 85.39           | 73.76        | 82.63        | 66.48        | <b>78.68</b> | 76.13         | 75.07        |
|                 | Gemma-2-9B | 0        | -              | 67.50        | 85.10           | 88.15        | 61.80        | 16.07        | 32.79        | —             | 37.90        |
| Rand            |            |          | 93.33          | 96.15        | 89.52           | 74.70        | 84.29        | 74.40        | 13.89        | 77.19         | 75.80        |
| TopK            |            |          | <b>94.47</b>   | 96.34        | <b>90.50</b>    | 75.19        | 84.25        | 74.50        | 61.14        | 74.82         | 75.24        |
| 4               |            | Div-S3   | 93.64          | 96.54        | 89.98           | 75.51        | 84.07        | <b>74.86</b> | 37.32        | <b>77.94</b>  | <b>76.13</b> |
|                 |            | Div      | 93.45          | 95.69        | 90.03           | 74.85        | <b>84.44</b> | 73.34        | 36.29        | 77.06         | 75.96        |
|                 |            | TopK-Div | 93.34          | <b>96.57</b> | 90.19           | <b>75.60</b> | 83.54        | 74.47        | <b>70.43</b> | 75.05         | 75.21        |
| 8               |            | Rand     | 93.30          | 96.09        | 89.39           | 75.98        | <b>84.34</b> | 74.48        | 24.36        | 79.23         | 76.28        |
|                 |            | TopK     | <b>94.20</b>   | 96.55        | <b>90.62</b>    | 76.14        | 83.57        | 75.36        | 71.00        | 77.59         | 75.55        |
|                 |            | Div-S3   | 93.38          | <b>96.60</b> | 90.10           | <b>76.98</b> | 83.56        | <b>75.62</b> | 45.86        | <b>79.79</b>  | <b>77.24</b> |
|                 |            | Div      | 93.41          | 95.94        | 89.90           | 76.60        | 84.22        | 74.69        | 42.07        | 79.05         | 76.65        |
|                 |            | TopK-Div | 94.04          | 96.58        | 90.48           | 76.53        | 83.85        | 75.16        | <b>76.32</b> | 77.64         | 76.24        |
| Mistral-7B-v0.3 |            | 0        | -              | 66.50        | 94.00           | 76.41        | 51.80        | 9.48         | 5.17         | —             | 30.50        |
|                 | Rand       |          | 91.00          | 94.02        | 82.77           | 69.83        | 48.78        | 37.20        | 12.14        | <b>76.70</b>  | 74.71        |
|                 | TopK       |          | <b>93.57</b>   | <b>96.17</b> | <b>85.21</b>    | 69.73        | 49.28        | 38.20        | 60.14        | 75.04         | 73.73        |
|                 | 4          | Div-S3   | 92.83          | 95.60        | 83.93           | <b>70.29</b> | <b>51.43</b> | 38.22        | 37.75        | 75.74         | 75.54        |
|                 |            | Div      | 91.98          | 94.15        | 82.98           | 70.15        | 49.49        | 37.50        | 34.89        | 75.96         | <b>75.83</b> |
|                 |            | TopK-Div | 92.73          | 95.90        | 84.55           | 69.91        | 49.99        | <b>38.45</b> | <b>71.46</b> | 74.43         | 73.16        |
|                 | 8          | Rand     | 92.49          | 95.35        | 83.69           | 71.65        | 47.86        | 36.32        | 22.18        | 77.30         | 75.54        |
|                 |            | TopK     | <b>93.61</b>   | <b>96.15</b> | <b>85.17</b>    | 71.88        | 48.43        | 37.35        | 70.50        | 77.05         | 75.44        |
|                 |            | Div-S3   | 92.79          | 96.10        | 84.38           | <b>72.47</b> | 48.57        | 36.60        | 45.86        | <b>77.71</b>  | <b>76.56</b> |
|                 |            | Div      | 92.55          | 95.10        | 84.27           | 72.04        | 48.33        | 36.12        | 39.14        | 77.67         | 76.30        |
|                 |            | TopK-Div | 93.47          | 96.11        | 84.85           | 71.81        | <b>48.60</b> | <b>37.81</b> | <b>77.93</b> | 77.44         | 75.22        |

Table 2: **(Comparison of different methods when demonstration and query come from different distribution)** Across tasks with  $K = 4$  shots, diversity-aware methods are more robust to out-of-distribution query. The performance drop from ID to OOD on `TopK` is in general larger than diversity-aware methods.

| Test         | Demo. | Rand     | TopK  | Div-S3       | Div          | TopK-Div     |              |
|--------------|-------|----------|-------|--------------|--------------|--------------|--------------|
| Llama-3.1-8B | SST-2 | SST-2    | 91.31 | <b>94.13</b> | 92.89        | 91.50        | 92.75        |
|              |       | IMDB     | 88.85 | <b>90.80</b> | 86.90        | 90.71        | <b>90.80</b> |
|              |       | Amazon   | 88.28 | 89.50        | <b>90.70</b> | 86.64        | 89.60        |
|              | CsQA  | CsQA     | 71.15 | 72.54        | 72.28        | 71.17        | <b>72.57</b> |
|              |       | ARC-Easy | 66.86 | 66.70        | 66.50        | 67.08        | <b>67.70</b> |
|              |       | SCIQ     | 74.00 | 72.70        | <b>74.61</b> | 74.47        | 73.87        |
|              | SQuAD | 72.11    | 71.40 | <b>73.67</b> | 72.79        | 71.60        |              |
| Gemma-2-9B   | SST-2 | SST-2    | 93.33 | <b>94.47</b> | 93.64        | 93.45        | 93.34        |
|              |       | IMDB     | 88.66 | 89.90        | 85.50        | 88.59        | <b>91.10</b> |
|              |       | Amazon   | 88.69 | 89.40        | 89.30        | <b>90.49</b> | 89.60        |
|              | CsQA  | CsQA     | 74.70 | 75.19        | 75.51        | 74.85        | <b>75.60</b> |
|              |       | ARC-Easy | 68.30 | 68.90        | 68.80        | 68.58        | <b>69.50</b> |
|              |       | SCIQ     | 75.80 | 75.24        | <b>76.13</b> | 75.96        | 75.21        |
|              | SQuAD | 73.63    | 73.60 | <b>76.07</b> | 74.64        | 73.50        |              |

using Amazon (a shopping review dataset) as demonstrations, `TopK-Div` surpasses `TopK`. A similar trend is observed in commonsense reason-

ing: replacing Commonsense QA (ID) demonstrations with ARC-Easy (OOD) increases the performance gap between `Div` and `TopK` from 0.4% to 1.0%. Text-to-SQL generation follows this pattern, with a larger improvement in OOD settings. Additionally, we note that `GSM-Plus-Mini` serves as an OOD setting for `GSM8K` (Math in Table 1), as they share the same training set. A larger improvement from adding diversity is also observed on `GSM-Plus-Mini`.

For reading comprehension, switching to an OOD demonstration dataset does not significantly widen the gap between `Div` and `TopK`, but `Div` still outperforms `TopK`. We provide additional out-of-distribution (OOD) results in Appendix C.4, which further reinforce our conclusions.

Beyond explicitly defined OOD settings, the contrast between the Amazon and SST-2 classification tasks in Table 1 further illustrates the impact of

Table 3: Relative improvement of TopK-Div over TopK on GeoQuery and SQuAD on different sets of the queries. For the GeoQuery dataset, we fine-tuned both base models on its training set. We categorized questions in testing set as “Easy” if the fine-tuned models correctly answered them in a 0-shot setting, and as “Hard” if these models failed to answer them correctly in the same 0-shot setting. We report the performance of both methods in a 4-shot setting. For SQuAD, we split the testing set only using the fine-tuned gemma-2-9B model, since fine-tuning the Llama-3.1-8B model yielded poor results. We observe that TopK-Div exhibits greater improvement on “Hard” examples.

| Split | Method   | Gemma-2-9B |       | Llama-3.1-8B |       |
|-------|----------|------------|-------|--------------|-------|
|       |          | GeoQuery   | SQuAD | GeoQuery     | SQuAD |
| Easy  | TopK     | 72.09      | 83.01 | 79.31        | 81.04 |
|       | TopK-Div | 77.91      | 82.66 | 83.71        | 79.65 |
| Hard  | TopK     | 56.29      | 20.00 | 51.52        | 24.44 |
|       | TopK-Div | 67.11      | 23.70 | 62.13        | 25.19 |

distributional differences on the effectiveness of diversity-based selection. While SST-2 consists of curated movie reviews with relatively homogeneous content—where TopK consistently outperforms diversity-based methods—Amazon reviews span heterogeneous domains such as electronics, books, and household items. This broader domain variability in Amazon leads to performance gains for diversity-driven methods like Rand, Div-S3 and Div, sometimes even surpassing TopK. These results are consistent with the patterns observed in Table 2 and provide additional empirical support for our Finding 2.

**Finding 3: Diversity performs better on harder examples.** Besides discussing the performance of diversity-aware methods (TopK-Div, Div, and even Rand) at task levels, we also analyze which specific examples benefit most from diversity. For this, we first provide a method to quantify the “difficulty level” of examples. Motivated by (Swayamdipta et al., 2020), we use whether a model can correctly answer a question after fine-tuning as an indicator of that question’s difficulty for a specific language model. Therefore, we fine-tuned the corresponding base model on the dataset’s training set using LoRA. Subsequently, based on whether this fine-tuned model could accurately answer questions in the testing set under a zero-shot setting, we classified these questions as “easy” or “hard”.

We examine this phenomenon in GEOQUERY and SQUAD, where TopK-Div consistently outperforms TopK. Table 3 shows diversity yields greater benefits on harder examples. In GEO-

QUERY, the absolute accuracy improvement of TopK-Div over TopK is 5.11% on easy examples (averaged across two models), increasing to 10.72% on hard examples. In SQUAD, while TopK-Div slightly underperforms TopK on easy examples by 0.87%, it outperforms TopK on hard examples by 2.23%.

### 3.2 Understanding the Role of Diversity: Beyond Coverage Effects

We investigate diversity’s impact on ICL by disentangling *coverage* from *mechanisms beyond coverage* at both example and task levels.

**Example-level analysis.** In the GeoQuery dataset, the observation that diversity performs better on harder examples (see Table 3) aligns with the notion of enhanced coverage: difficult examples often require modeling more nuanced or rare local structures (Levy et al., 2023; Gupta et al., 2023), which diversity-based methods are more likely to capture.

However, in SQuAD, we observe a different pattern. Even when  $k = 1$ , TopK underperforms compared to Rand and Div, suggesting that coverage alone is insufficient to explain the performance gap. To probe this, we remove irrelevant noisy examples from the SQuAD dataset and rerun the comparison. This cleaning significantly improves TopK and TopK-Div, but has minimal impact on Rand and Div—indicating that the strength of diversity-based methods extends beyond simple structural alignment with the query.

We present detailed results and analysis in Appendix C.5 to support this claim.

**Task-level analysis.** As shown in Table 7, increasing the number of demonstrations (e.g., from 4 to 8 shots) magnifies the benefits of diversity. While coverage-based reasoning suggests this may be due to broader inclusion of features, our findings point to a richer effect.

Specifically, when given more demonstrations, models appear to better synthesize the overall conceptual structure of the task. Diversity enables this by exposing the model to varied facets of the task distribution, which helps form a more general and transferable representation. In contrast, with fewer shots, the model has limited capacity to form such abstractions, and similarity alone may suffice.

We further justify this effect theoretically in Appendix E, showing that diversity supports a form of generalization that cannot be fully explained by coverage alone.

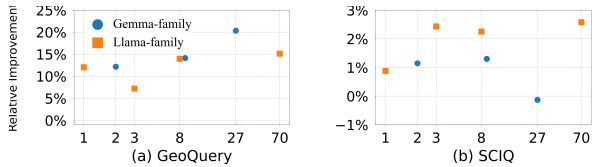


Figure 3: The relative improvement of diversity-aware methods over TopK. **Left:** relative improvement of TopK-Div over TopK on GeoQuery standard split. **Right:** relative improvement of Div over TopK on SCIQ.

### 3.3 Practical insights

In summary, our findings provide the following guidelines for selecting demonstrations in ICL:

- Leverage similarity for simple tasks.** When the task is relatively easy (e.g., sentiment classification) and the model already exhibits sufficient ability, selecting demonstrations purely based on similarity to the query is generally sufficient to elicit strong performance.
- Use diversity to bridge distribution gaps.** When there is a significant mismatch between the test distribution and the available demonstration pool, incorporating diversity in selection helps the model generalize better by exposing it to a broader range of examples.
- Favor diversity for complex or knowledge-intensive tasks.** For tasks that require the model to extract and apply task-solving knowledge (e.g., math or reading), selecting diverse demonstrations provides broader coverage of relevant patterns or skills.
- Adapt to noise levels.** For *low-noise datasets*, coverage-oriented (similarity-based) selection is more effective, as it aligns demonstrations closely with the query and helps the model lock onto the correct input-output mapping. In contrast, for *high-noise datasets*, increased diversity is beneficial to reduce overfitting to spurious correlations and enhance robustness.

### 3.4 Ablation studies

We examine how diversity-aware gains over TopK vary with LLM scale, as larger models might be less sensitive to data selection, potentially diminishing diversity’s advantages.

Figure 3 shows the relative improvement on GeoQuery standard split and SCIQ, where we observe the clear benefit of the diversity-aware method in Section 3.1, on Llama-3.1/3.2 and Gemma-2 families with different sizes. We observe that in these two tasks, in general, the relative improvement does

not decrease that much even if the model scales up, which indicates the importance of understanding the role of diversity in demonstration selection.

In Appendix D, we present experiments on additional models. We also report results across a range of settings, including fine-grained variations in  $k$ , different subset sizes for the Div method, fixed training sets, and changes in the embedding or decoding strategies. In particular, we implement a **purely** diversity-based method, K-Means, whose diversity score can exceed that of Div. Its superior performance on Math and Reading tasks further supports Finding 1. These ablation studies consistently reinforce our main findings, demonstrating the generality and robustness of our conclusions.

## 4 Conclusion, limitations, and future works

We investigate the role of diversity in retrieval-based demonstration selection for ICL. Across a wide range of tasks and model families, we find incorporating diversity into selection strategies consistently improves performance, especially when the task is difficult, the query is challenging, or there is a distribution shift between the query and available demonstrations. These findings are further supported by comprehensive ablation studies.

In addition, we provide theoretical justification for when and why diversity offers advantages over purely similarity-based selection. Together, our empirical and theoretical insights offer practical guidance for selecting effective demonstrations in ICL and deepen the understanding of diversity’s role in prompting large language models.

Note that the internal mechanism behind why diversity benefits still remains unclear. Part of our findings can be explained by coverage, which is aligned with previous literature, but the superior performance on math, reading comprehension, and OOD generalization, cannot be explained by simply incentivizing coverage. Potential future research directions include both theoretical and empirical explorations into why diversity aids demonstration selection beyond coverage. This could involve deeper analysis of model representations, interactions between diverse demonstrations, or alternative explanations grounded in information theory or representation learning. Additionally, our diversity heuristic is tested on English text only; cross-lingual robustness is left for future work.

## 643 Limitations

644 While our results consistently show that diversity-  
645 aware selection improves performance on challeng-  
646 ing tasks and under distribution shift, the mech-  
647 anism remains incompletely understood: cover-  
648 age can explain part of the gains, yet the im-  
649 provements on math, reading comprehension, and  
650 OOD settings are not fully attributable to coverage  
651 alone. Establishing a more causal account (e.g.,  
652 via representation-level analyses or controlled syn-  
653 thetic setups) is left for future work.

654 Our empirical study focuses on English bench-  
655 marks; cross-lingual robustness of the proposed  
656 heuristics remains to be explored.

657 Beyond these, our theoretical analysis relies on a  
658 linear model to provide an interpretable account of  
659 why diversity can help. However, this abstraction  
660 omits many properties of modern LLMs (nonlinear-  
661 ity, depth, attention, and optimization dynamics),  
662 so the theory should be viewed as offering qualita-  
663 tive insights rather than quantitative predictions.

## 664 Ethics Statement

665 This research fully aligns with the ethical princi-  
666 ples of ACL. Our work systematically investigates  
667 the role of sample-selection diversity in *In-context*  
668 *learning (ICL)*, aiming to improve the performance,  
669 robustness, and reliability of large language mod-  
670 els.

671 Our work studies diversity-aware demonstration  
672 selection for in-context learning as a foundational  
673 method, rather than proposing or releasing a de-  
674 ployable system. As such, it does not directly intro-  
675 duce new capabilities for disinformation, surveil-  
676 lance, or privacy attacks; we do not release new  
677 models, collect user data, or construct datasets con-  
678 taining personal information.

679 The primary motivation is to contribute posi-  
680 tively to society and human well-being. By demon-  
681 strating that diversity-aware selection of in-context  
682 examples can lead to improvements on complex  
683 tasks (such as mathematical reasoning and code  
684 generation) and out-of-distribution queries, we  
685 hope to foster AI systems that generalize better,  
686 thereby serving societal applications in research,  
687 education, software development, and beyond. We  
688 pay particular attention to underrepresented or chal-  
689 lenging settings, in line with the ACL principle of  
690 giving emphasis to less-advantaged groups.

691 In striving for scientific excellence, we adhere  
692 to methodological rigor, transparency, and repro-

693 ducibility. Our conclusions are supported by sys-  
694 tematic experiments across multiple tasks, datasets,  
695 and models (including LLaMA 3.1, Gemma 2,  
696 Mistral-v0.3), and by a theoretical framework that  
697 elucidates why diversity helps. We use publicly  
698 available benchmark datasets (e.g. SST-2, GSM8K,  
699 GeoQuery) and open-source models, and we pro-  
700 vide full details of experimental design, hyperpa-  
701 rameters, and evaluation procedures in the paper  
702 and appendix.

703 We also commit to fairness and non-  
704 discrimination. Although bias mitigation is  
705 not a direct focus, our findings suggest that  
706 diversity-aware in-context selection can improve  
707 out-of-distribution robustness, potentially helping  
708 models maintain performance even for underrepre-  
709 sented groups or non-mainstream distributions. We  
710 view this as a positive step toward more equitable  
711 AI systems.

712 Privacy and respect for intellectual labor are also  
713 core commitments. Our study uses only publicly  
714 available, anonymized datasets; no new personal  
715 or sensitive data were collected, and no human  
716 subjects were involved. We fully cite all utilized  
717 datasets, models, and prior work, giving due credit  
718 to others' contributions.

719 In summary, we believe this work is a respon-  
720 sible and beneficial contribution toward building  
721 more robust and trustworthy large language models.  
722 We have carefully considered the relevant ethical  
723 dimensions and commit to conducting our research  
724 according to the scientific and ethical standards  
725 expected under ACL's Code of Ethics.

## References

- 726  
727 Ekin Akyürek, Dale Schuurmans, Jacob Andreas,  
728 Tengyu Ma, and Denny Zhou. 2023. **What learn-  
729 ing algorithm is in-context learning? investigations  
730 with linear models.** In *The Eleventh International  
731 Conference on Learning Representations, ICLR 2023,  
732 Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- 733 Shengnan An, Zeqi Lin, Qiang Fu, Bei Chen, Nanning  
734 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), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 11027–11052. Association for Computational Linguistics.
- 741 Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen,  
742 Nanning Zheng, Weizhu Chen, and Jian-Guang Lou.  
743 2023b. **Skill-based few-shot selection for in-context  
744 learning.** In *Proceedings of the 2023 Conference*

|     |  |  |     |
|-----|--|--|-----|
| 745 | <i>on Empirical Methods in Natural Language Processing</i> , pages 13472–13492, Singapore. Association for Computational Linguistics.  | gradient descent as meta-optimizers. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 4005–4019.  | 801 |
| 746 |  |  | 802 |
| 747 |  |  | 803 |
| 748 | Sanjeev Arora and Anirudh Goyal. 2023. A theory for emergence of complex skills in language models. <i>arXiv preprint arXiv:2307.15936</i> .   | Gilad Deutch, Nadav Magar, Tomer Natan, and Guy Dar. 2024. In-context learning and gradient descent revisited. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1017–1028. | 804 |
| 749 |  |  | 805 |
| 750 |  |  | 806 |
| 751 | Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.  | Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke, Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo Rezende, Yoshua Bengio, Michael Mozer, and Sanjeev Arora. 2024. Metacognitive capabilities of llms: An exploration in mathematical problem solving. <i>arXiv preprint arXiv:2405.12205</i> .                  | 807 |
| 752 |  |  | 808 |
| 753 |  |  | 809 |
| 754 |  |  | 810 |
| 755 |  |  | 811 |
| 756 |  |  | 812 |
| 757 | Ting-Yun Chang and Robin Jia. 2023. Data curation alone can stabilize in-context learning. In <i>The 61st Annual Meeting Of The Association For Computational Linguistics</i> .  | Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .   | 813 |
| 758 |  |  | 814 |
| 759 |  |  | 815 |
| 760 |  |  | 816 |
| 761 | Laming Chen, Guoxin Zhang, and Eric Zhou. 2018. Fast greedy map inference for determinantal point process to improve recommendation diversity. <i>Advances in Neural Information Processing Systems</i> , 31.  | Shuzheng Gao, Xin-Cheng Wen, Cuiyun Gao, Wenxuan Wang, and Michael R Lyu. 2023. Constructing effective in-context demonstration for code intelligence tasks: An empirical study. <i>CoRR</i> .   | 817 |
| 762 |  |  | 818 |
| 763 |  |  | 819 |
| 764 |  |  | 820 |
| 765 |  |  | 821 |
| 766 | Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. <i>arXiv preprint arXiv:2107.03374</i> .  | Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn in-context? a case study of simple function classes. <i>Advances in Neural Information Processing Systems</i> , 35:30583–30598.   | 822 |
| 767 |  |  | 823 |
| 768 |  |  | 824 |
| 769 |  |  | 825 |
| 770 |  |  | 826 |
| 771 |  |  | 827 |
| 772 | Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. 2023. On the relation between sensitivity and accuracy in in-context learning. In <i>2023 Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 155–167. Association for Computational Linguistics (ACL).   | Petros Stylianos Giouroukis, Alexios Gidiotis, and Grigorios Tsoumakas. 2025. <b>Dual: Diversity and uncertainty active learning for text summarization</b> . <i>Preprint</i> , arXiv:2503.00867.  | 828 |
| 773 |  |  | 829 |
| 774 |  |  | 830 |
| 775 |  |  | 831 |
| 776 |  |  | 832 |
| 777 |  |  | 833 |
| 778 | Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv preprint arXiv:1803.05457</i> .  | Shivanshu Gupta, Matt Gardner, and Sameer Singh. 2023. Coverage-based example selection for in-context learning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 13924–13950.  | 834 |
| 779 |  |  | 835 |
| 780 |  |  | 836 |
| 781 |  |  | 837 |
| 782 |  |  | 838 |
| 783 | 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. <i>arXiv preprint arXiv:2110.14168</i> .  | Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, and 1 others. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .  | 839 |
| 784 |  |  | 840 |
| 785 |  |  | 841 |
| 786 |  |  | 842 |
| 787 |  |  | 843 |
| 788 |  |  | 844 |
| 789 | Ryan A. Cook, John P. Lalor, and Ahmed Abbasi. 2025. <b>No simple answer to data complexity: An examination of instance-level complexity metrics for classification tasks</b> . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 2553–2573, Albuquerque, New Mexico. Association for Computational Linguistics. | Hui Jiang. 2023. A latent space theory for emergent abilities in large language models. <i>arXiv preprint arXiv:2304.09960</i> .   | 845 |
| 790 |  |  | 846 |
| 791 |  |  | 847 |
| 792 |  |  | 848 |
| 793 |  |  | 849 |
| 794 |  |  | 850 |
| 795 |  |  | 851 |
| 796 |  |  | 852 |
| 797 |  |  | 853 |
| 798 | Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2023. Why can gpt learn in-context? language models secretly perform   | Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781.    | 854 |
| 799 |  |  | 855 |
| 800 |  |  | 856 |

|     |   |     |
|-----|---|-----|
| 854 | Lilly Kumari, Shengjie Wang, Arnav Das, Tianyi Zhou, and Jeff Bilmes. 2024. <b>An end-to-end submodular framework for data-efficient in-context learning</b> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3293–3308, Mexico City, Mexico. Association for Computational Linguistics.   | 911 |
| 855 |   | 912 |
| 856 |   | 913 |
| 857 |   | 914 |
| 858 |   | 915 |
| 859 |   |     |
| 860 |   |     |
| 861 | Itay Levy, Ben Bogin, and Jonathan Berant. 2023. Diverse demonstrations improve in-context compositional generalization. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1401–1422.   | 916 |
| 862 |   | 917 |
| 863 |   | 918 |
| 864 |   | 919 |
| 865 |   | 920 |
| 866 |   | 921 |
| 867 | Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024. <b>Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers</b> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 2961–2984. Association for Computational Linguistics. | 922 |
| 868 |   | 923 |
| 869 |   | 924 |
| 870 |   | 925 |
| 871 |   | 926 |
| 872 |   | 927 |
| 873 |   | 928 |
| 874 |   | 929 |
| 875 |   | 930 |
| 876 | Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. 2023a. The closeness of in-context learning and weight shifting for softmax regression. <i>arXiv preprint arXiv:2304.13276</i> .   | 931 |
| 877 |   | 932 |
| 878 |   | 933 |
| 879 |   | 934 |
| 880 |   | 935 |
| 881 | Xiaonan Li and Xipeng Qiu. 2023a. Finding support examples for in-context learning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 6219–6235.  | 936 |
| 882 |   | 937 |
| 883 |   | 938 |
| 884 | Xiaonan Li and Xipeng Qiu. 2023b. <b>Finding support examples for in-context learning</b> . <i>Preprint</i> , arXiv:2302.13539.   | 939 |
| 885 |   | 940 |
| 886 |   | 941 |
| 887 | Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Pappaliopoulos, and Samet Oymak. 2023b. Transformers as algorithms: Generalization and stability in in-context learning. In <i>International Conference on Machine Learning</i> , pages 19565–19594. PMLR.   | 942 |
| 888 |   | 943 |
| 889 |   | 944 |
| 890 |   | 945 |
| 891 |   | 946 |
| 892 | Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. <i>arXiv preprint arXiv:2305.20050</i> .  | 947 |
| 893 |   | 948 |
| 894 |   | 949 |
| 895 |   | 950 |
| 896 |   | 951 |
| 897 | Hui Lin and Jeff Bilmes. 2011. A class of submodular functions for document summarization. In <i>Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies</i> , pages 510–520.  | 952 |
| 898 |   | 953 |
| 899 |   | 954 |
| 900 |   | 955 |
| 901 |   | 956 |
| 902 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? <i>arXiv preprint arXiv:2101.06804</i> .   | 957 |
| 903 |   | 958 |
| 904 |   | 959 |
| 905 |   | 960 |
| 906 |   | 961 |
| 907 | Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.   | 962 |
| 908 |   | 963 |
| 909 |   | 964 |
| 910 |   | 965 |
|     |   | 966 |
|     | Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. <i>arXiv preprint arXiv:2104.08786</i> .   |     |
|     | 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.   |     |
|     | Man Luo, Xin Xu, Zhuyun Dai, Panupong Pasupat, Mehran Kazemi, Chitta Baral, Vaiva Imbrasaite, and Vincent Y Zhao. 2023. Dr. icl: Demonstration-retrieved in-context learning. <i>arXiv preprint arXiv:2305.14128</i> .  |     |
|     | Man Luo, Xin Xu, Yue Liu, Panupong Pasupat, and Mehran Kazemi. 2024. <b>In-context learning with retrieved demonstrations for language models: A survey</b> . <i>Preprint</i> , arXiv:2401.11624.   |     |
|     | Xuetao Ma, Wenbin Jiang, and Hua Huang. 2025. <b>Problem-solving logic guided curriculum in-context learning for llms complex reasoning</b> . <i>Preprint</i> , arXiv:2502.15401.   |     |
|     | Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies</i> , pages 142–150.   |     |
|     | Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In <i>Proceedings of the 7th ACM conference on Recommender systems</i> , pages 165–172.   |     |
|     | Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064.   |     |
|     | Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. 2023. What in-context learning “learns” in-context: Disentangling task recognition and task learning. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8298–8319.   |     |
|     | Adarsh Prasad, Stefanie Jegelka, and Dhruv Batra. 2014. Submodular meets structured: Finding diverse subsets in exponentially-large structured item sets. <i>Advances in Neural Information Processing Systems</i> , 27.  |     |
|     | Chengwei Qin, Aston Zhang, Chen Chen, Anirudh Dagar, and Wenming Ye. 2023. In-context learning with iterative demonstration selection. <i>arXiv preprint arXiv:2310.09881</i> .   |     |

|      |   |      |
|------|---|------|
| 967  | Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <b>SQuAD: 100,000+ questions for machine comprehension of text</b> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.   | 1022 |
| 968  |   | 1023 |
| 969  |   | 1024 |
| 970  |   | 1025 |
| 971  |   | 1026 |
| 972  |   | 1027 |
| 973  | N Reimers. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .   | 1028 |
| 974  |   | 1029 |
| 975  |   | 1030 |
| 976  | Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2655–2671.   | 1031 |
| 977  |   | 1032 |
| 978  |   | 1033 |
| 979  |   | 1034 |
| 980  |   | 1035 |
| 981  |   | 1036 |
| 982  | Alexander Scarlatos and Andrew Lan. 2023. Reticl: Sequential retrieval of in-context examples with reinforcement learning. <i>arXiv preprint arXiv:2305.14502</i> .   | 1037 |
| 983  |   | 1038 |
| 984  |   | 1039 |
| 985  |   | 1040 |
| 986  | Ozan Sener and Silvio Savarese. 2018. <b>Active learning for convolutional neural networks: A core-set approach</b> . <i>Preprint</i> , arXiv:1708.00489.   | 1041 |
| 987  |   | 1042 |
| 988  |   | 1043 |
| 989  | Lingfeng Shen, Aayush Mishra, and Daniel Khashabi. 2023. Do pretrained transformers really learn in-context by gradient descent? <i>arXiv preprint arXiv:2310.08540</i> .   | 1044 |
| 990  |   | 1045 |
| 991  |   | 1046 |
| 992  |   | 1047 |
| 993  | Lei Shi and Yi-Dong Shen. 2016. Diversifying convex transductive experimental design for active learning. In <i>IJCAI</i> , pages 1997–2003.  | 1048 |
| 994  |   | 1049 |
| 995  |   | 1050 |
| 996  | Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, and 432 others. 2023. <b>Beyond the imitation game: Quantifying and extrapolating the capabilities of language models</b> . <i>Preprint</i> , arXiv:2206.04615. | 1051 |
| 997  |   | 1052 |
| 998  |   | 1053 |
| 999  |   | 1054 |
| 1000 |   | 1055 |
| 1001 |   | 1056 |
| 1002 |   | 1057 |
| 1003 |   | 1058 |
| 1004 |   | 1059 |
| 1005 |   | 1060 |
| 1006 | Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. <b>Selective annotation makes language models better few-shot learners</b> . In <i>The Eleventh International Conference on Learning Representations</i> .  | 1061 |
| 1007 |   | 1062 |
| 1008 |   | 1063 |
| 1009 |   | 1064 |
| 1010 |   | 1065 |
| 1011 |   | 1066 |
| 1012 | Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A. Smith, and Yejin Choi. 2020. <b>Dataset cartography: Mapping and diagnosing datasets with training dynamics</b> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 9275–9293. Association for Computational Linguistics.  | 1067 |
| 1013 |   | 1068 |
| 1014 |   | 1069 |
| 1015 |   | 1070 |
| 1016 |   | 1071 |
| 1017 |   | 1072 |
| 1018 |   | 1073 |
| 1019 |   | 1074 |
| 1020 | Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question   | 1075 |
| 1021 |   | 1076 |
|      | answering challenge targeting commonsense knowledge. In <i>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)</i> , pages 4149–4158.  | 1022 |
|      |   | 1023 |
|      |   | 1024 |
|      |   | 1025 |
|      |   | 1026 |
|      |   | 1027 |
|      | Lappoon R Tang and Raymond J Mooney. 2001. Using multiple clause constructors in inductive logic programming for semantic parsing. In <i>European Conference on Machine Learning</i> , pages 466–477. Springer.   | 1028 |
|      |   | 1029 |
|      |   | 1030 |
|      |   | 1031 |
|      |   | 1032 |
|      | Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .  | 1033 |
|      |   | 1034 |
|      |   | 1035 |
|      |   | 1036 |
|      |   | 1037 |
|      |   | 1038 |
|      | Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023a. Transformers learn in-context by gradient descent. In <i>International Conference on Machine Learning</i> , pages 35151–35174. PMLR.   | 1039 |
|      |   | 1040 |
|      |   | 1041 |
|      |   | 1042 |
|      |   | 1043 |
|      |   | 1044 |
|      | Johannes Von Oswald, Maximilian Schlegel, Alexander Meulemans, Seijin Kobayashi, Eyvind Niklasson, Nicolas Zucchet, Nino Scherrer, Nolan Miller, Mark Sandler, Max Vladymyrov, and 1 others. 2023b. Uncovering mesa-optimization algorithms in transformers. <i>arXiv preprint arXiv:2309.05858</i> .   | 1045 |
|      |   | 1046 |
|      |   | 1047 |
|      |   | 1048 |
|      |   | 1049 |
|      |   | 1050 |
|      | Zhijing Wan, Zhixiang Wang, Yuran Wang, Zheng Wang, Hongyuan Zhu, and Shin’ichi Satoh. 2024. <b>Contributing dimension structure of deep feature for coresets selection</b> . <i>Preprint</i> , arXiv:2401.16193.   | 1051 |
|      |   | 1052 |
|      |   | 1053 |
|      |   | 1054 |
|      | Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. 2024. <b>Diversity measurement and subset selection for instruction tuning datasets</b> . <i>Preprint</i> , arXiv:2402.02318.  | 1055 |
|      |   | 1056 |
|      |   | 1057 |
|      |   | 1058 |
|      |   | 1059 |
|      | Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. In <i>Workshop on Efficient Systems for Foundation Models@ ICML2023</i> .  | 1060 |
|      |   | 1061 |
|      |   | 1062 |
|      |   | 1063 |
|      |   | 1064 |
|      |   | 1065 |
|      | Xinyi Wang, Wanrong Zhu, and William Yang Wang. 2023. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. <i>arXiv preprint arXiv:2301.11916</i> , page 3.   | 1066 |
|      |   | 1067 |
|      |   | 1068 |
|      |   | 1069 |
|      |   | 1070 |
|      | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.   | 1071 |
|      |   | 1072 |
|      |   | 1073 |
|      |   | 1074 |
|      |   | 1075 |
|      |   | 1076 |

|      |  |  |      |
|------|--|--|------|
| 1077 | Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017.                | Suqin Yuan, Lei Feng, Bo Han, and Tongliang Liu.                     | 1133 |
| 1078 | Crowdsourcing multiple choice science questions.                     | 2025. <b>Enhancing sample selection against label</b>                | 1134 |
| 1079 | In <i>Proceedings of the 3rd Workshop on Noisy User-</i>             | <b>noise by cutting mislabeled easy examples</b> . <i>Preprint</i> , | 1135 |
| 1080 | <i>generated Text</i> , pages 94–106.                                | arXiv:2502.08227.  | 1136 |
| 1081 | Noam Wies, Yoav Levine, and Amnon Shashua. 2023.                     | John M Zelle and Raymond J Mooney. 1996. Learning                    | 1137 |
| 1082 | The learnability of in-context learning. <i>Advances in</i>          | to parse database queries using inductive logic pro-                 | 1138 |
| 1083 | <i>Neural Information Processing Systems</i> , 36:36637–             | gramming. In <i>Proceedings of the national conference</i>           | 1139 |
| 1084 | 36651.   | <i>on artificial intelligence</i> , pages 1050–1055.                 | 1140 |
| 1085 | Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Ling-                    | Donglin Zhan, Leonardo F. Toso, and James Ander-                     | 1141 |
| 1086 | peng Kong. 2023. Self-adaptive in-context learn-                     | son. 2025. <b>Coreset-based task selection for sample-</b>           | 1142 |
| 1087 | ing: An information compression perspective for in-                  | <b>efficient meta-reinforcement learning</b> . <i>Preprint</i> ,     | 1143 |
| 1088 | context example selection and ordering. In <i>Proceed-</i>           | arXiv:2502.02332.  | 1144 |
| 1089 | <i>ings of the 61st Annual Meeting of the Association for</i>        | Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and                      | 1145 |
| 1090 | <i>Computational Linguistics (Volume 1: Long Papers)</i> ,           | Zhaoran Wang. 2023a. What and how does in-                           | 1146 |
| 1091 | pages 1423–1436.   | context learning learn? bayesian model averaging,                    | 1147 |
| 1092 | Sang Michael Xie, Aditi Raghunathan, Percy Liang,                    | parameterization, and generalization. <i>arXiv preprint</i>          | 1148 |
| 1093 | and Tengyu Ma. 2022. <b>An explanation of in-context</b>             | <i>arXiv:2305.19420</i> .  | 1149 |
| 1094 | <b>learning as implicit bayesian inference</b> . In <i>The Tenth</i> | Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex                        | 1150 |
| 1095 | <i>International Conference on Learning Representa-</i>              | Smola. 2023b. <b>Automatic chain of thought prompt-</b>              | 1151 |
| 1096 | <i>tions, ICLR 2022, Virtual Event, April 25-29, 2022</i> .          | <b>ing in large language models</b> . In <i>The Eleventh Inter-</i>  | 1152 |
| 1097 | OpenReview.net.  | <i>national Conference on Learning Representations,</i>              | 1153 |
| 1098 | Zhao Yang, Yuanzhe Zhang, Dianbo Sui, Cao Liu, Jun                   | <i>ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . Open-              | 1154 |
| 1099 | Zhao, and Kang Liu. 2023. Representative demon-                      | Review.net.  | 1155 |
| 1100 | stration selection for in-context learning with two-                 | Haoyu Zhao, Simran Kaur, Dingli Yu, Anirudh Goyal,                   | 1156 |
| 1101 | stage determinantal point process. In <i>Proceedings</i>             | and Sanjeev Arora. 2025. Can models learn skill                      | 1157 |
| 1102 | <i>of the 2023 Conference on Empirical Methods in</i>                | composition from examples? <i>Advances in Neural</i>                 | 1158 |
| 1103 | <i>Natural Language Processing</i> , pages 5443–5456.                | <i>Information Processing Systems</i> , 37:102393–102427.            | 1159 |
| 1104 | Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu,                  | Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and                   | 1160 |
| 1105 | Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. <b>An</b>              | Sameer Singh. 2021. Calibrate before use: Improv-                    | 1161 |
| 1106 | <b>empirical study of GPT-3 for few-shot knowledge-</b>              | ing few-shot performance of language models. In                      | 1162 |
| 1107 | <b>based VQA</b> . In <i>Thirty-Sixth AAAI Conference on</i>         | <i>International conference on machine learning</i> , pages          | 1163 |
| 1108 | <i>Artificial Intelligence, AAAI 2022, Thirty-Fourth</i>             | 12697–12706. PMLR.   | 1164 |
| 1109 | <i>Conference on Innovative Applications of Artificial In-</i>       |  |      |
| 1110 | <i>telligence, IAAI 2022, The Twelveth Symposium on</i>              |  |      |
| 1111 | <i>Educational Advances in Artificial Intelligence, EAAI</i>         |  |      |
| 1112 | <i>2022 Virtual Event, February 22 - March 1, 2022</i> ,             |  |      |
| 1113 | pages 3081–3089. AAAI Press.   |  |      |
| 1114 | Xi Ye, Srinivasan Iyer, Asli Celikyilmaz, Ves Stoy-                  |  |      |
| 1115 | anov, Greg Durrett, and Ramakanth Pasunuru. 2023.                    |  |      |
| 1116 | Complementary explanations for effective in-context                  |  |      |
| 1117 | learning. <i>Findings of the Association for Computa-</i>            |  |      |
| 1118 | <i>tional Linguistics: ACL 2023</i> .                                |  |      |
| 1119 | Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyun-                    |  |      |
| 1120 | soo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee,                     |  |      |
| 1121 | and Taek Kim. 2022. Ground-truth labels matter: A                    |  |      |
| 1122 | deeper look into input-label demonstrations. In <i>Pro-</i>          |  |      |
| 1123 | <i>ceedings of the 2022 Conference on Empirical Meth-</i>            |  |      |
| 1124 | <i>ods in Natural Language Processing</i> , pages 2422–              |  |      |
| 1125 | 2437.  |  |      |
| 1126 | Dingli Yu, Simran Kaur, Arushi Gupta, Jonah Brown-                   |  |      |
| 1127 | Cohen, Anirudh Goyal, and Sanjeev Arora. 2024.                       |  |      |
| 1128 | <b>SKILL-MIX: a flexible and expandable family of</b>                |  |      |
| 1129 | <b>evaluations for AI models</b> . In <i>The Twelfth Inter-</i>      |  |      |
| 1130 | <i>national Conference on Learning Representations,</i>              |  |      |
| 1131 | <i>ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . Open-            |  |      |
| 1132 | Review.net.  |  |      |

1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197

**Contents**

- 1 Introduction 1**
- 1.1 Our contributions . . . . . 2
- 2 Background and notations 3**
- 3 Experiments and findings 4**
- 3.1 Main findings . . . . . 5
- 3.2 Understanding the Role of Diversity: Beyond Coverage Effects . . . . . 7
- 3.3 Practical insights . . . . . 8
- 3.4 Ablation studies . . . . . 8
- 4 Conclusion, limitations, and future works 8**
- A More related works 15**
- B More experiment details 15**
- B.1 Prompt template . . . . . 15
- B.2 Dataset details . . . . . 15
- B.3 Evaluation details . . . . . 15
- C Additional experiments 17**
- C.1 Results of 0/1-shot . . . . . 17
- C.2 Additional supplement of Finding 1: Adjusting subset sizes for Div . . . . . 18
- C.3 Additional supplement of Finding 1: Adjusting  $\alpha$  for TopK-Div . . . . . 18
- C.4 More results on OOD setting . . . . . 18
- C.5 Results on perturbation of datasets . . . . . 19
- D Additional ablation studies 20**
- D.1 Results on more models . . . . . 20
- D.2 Changing the number of shots . . . . . 21
- D.3 More diversity-aware method . . . . . 24
- D.4 Abalations on the size of training set . . . . . 24
- D.5 Ablations on “better” embeddings . . . . . 24
- D.6 Decoding method . . . . . 25
- E Theoretical justification and simulations 26**
- E.1 Proof of Theorem E.2: justification example I . . . . . 27
- E.2 Proof of Theorem E.3: justification example II . . . . . 29
- E.3 Experiment settings . . . . . 31
- E.4 Result and discussions . . . . . 32

|   |  |
|---|--|
| <b>A More related works</b>   | 1198   |
| <b>Demonstration selection</b> Retrieval-based demonstration selection for ICL has long been studied, and the most notable methods are the <i>similarity</i> -based methods (Liu et al., 2021; Yang et al., 2022; Wu et al., 2023; Qin et al., 2023). These are often augmented by trainable deep learning retrievers aimed at capturing core skills or features beyond mere semantic similarity (Karpukhin et al., 2020; Rubin et al., 2022; Luo et al., 2023; Scarlatos and Lan, 2023; An et al., 2023b), or by incorporating LLM feedback for refinement (Li and Qiu, 2023a; Chen et al., 2023; Wang et al., 2023). Conversely, diversity-based, or more accurately, coverage-based methods are less prevalent in retrieval-based selection. Existing studies in this vein typically address tasks with clear local structures where feature coverage is advantageous (Levy et al., 2023; Ye et al., 2023; Gupta et al., 2023; An et al., 2023a). For non-retrieval-based ICL, where a fixed set of demonstrations is selected for a specific task, diversity is recognized as beneficial (Zhang et al., 2023b; Gao et al., 2023; Su et al., 2023; Yang et al., 2023). | 1199<br>1200<br>1201<br>1202<br>1203<br>1204<br>1205<br>1206<br>1207<br>1208<br>1209 |
| <b>Understanding in-context learning</b> Efforts to understand ICL span both theoretical and empirical investigations. Theoretical perspectives often frame ICL as either a Bayesian inference procedure (Xie et al., 2022; Wang et al.; Wies et al., 2023; Jiang, 2023; Zhang et al., 2023a) or an implicit form of meta-optimization akin to gradient descent (Dai et al., 2023; Von Oswald et al., 2023a,b; Deutch et al., 2024; Shen et al., 2023). Research on ICL for regression tasks (Garg et al., 2022; Li et al., 2023b,a; Akyürek et al., 2023) provides valuable insights; notably, (Akyürek et al., 2023) suggest transformers can identify min-norm solutions in-context for linear regression, a finding that supports the role of demonstration diversity. Empirical studies have further examined factors such as input-label mapping (Min et al., 2022; Yoo et al., 2022; Pan et al., 2023), the influence of demonstration order (Lu et al., 2022; Liu et al., 2024), and the importance of calibration for ICL efficacy (Zhao et al., 2021).  | 1210<br>1211<br>1212<br>1213<br>1214<br>1215<br>1216<br>1217<br>1218<br>1219         |
| <b>B More experiment details</b>  | 1220   |
| <b>B.1 Prompt template</b>  | 1221   |
| Table 4 lists the template we use for different tasks. We take $K = 2$ as an example.   | 1222   |
| <b>B.2 Dataset details</b>  | 1223   |
| To reduce computational cost, we performed random sampling on both the <i>demo</i> and <i>test</i> set for classification, multiple-choice and reading tasks. For classification tasks, the sampled datasets from IMDB and SST-2 are consistent with (Chang and Jia, 2023). A fixed random seed of 42 was used for all sampling procedures. For math tasks, since the <i>test</i> set sizes of PRM800K and GSM8K datasets are close to the sampled <i>test</i> set sizes of other tasks, we directly used their existing <i>demo</i> and <i>test</i> set. Detailed sampling statistics are provided in Table 5.   | 1224<br>1225<br>1226<br>1227<br>1228<br>1229   |
| <b>B.3 Evaluation details</b>   | 1230   |
| For the sentiment classification task (classification), given the prompt listed in Table 4, we compute the logit for “great” and “terrible” respectively, and predict the sentiment to be positive if the logit for “great” is larger than that for “terrible”, and vice versa. We report the accuracy metric.  | 1231<br>1232<br>1233   |
| For commonsense reasoning tasks (multiple-choice), given the prompt, we compute the average cross-entropy loss on each given option, conditioned on the prompt. Then we pick the option with the smallest average cross-entropy loss. We report the accuracy metric.  | 1234<br>1235<br>1236   |
| For reading comprehension (generation), given the prompt, we generate the answer using greedy decoding. We stop if we generate one of the following string: “\n\n”, “\n\n\n”, "Support", "Support:", "Question", "Question:". We compare the generated answer with the gold answer, and report the exact match metric. There are several optional answers for the squad test sample, if the generated answer exactly matches one of them, we consider it correct.   | 1237<br>1238<br>1239<br>1240<br>1241   |
| For text to SQL (generation), given the prompt, we generate the answer using greedy decoding. We stop if we generate one of the following string: “\n\n”, “\n\n\n”, "Question", "Question:". We compare the generated answer with the gold answer, and report the exact match metric.   | 1242<br>1243<br>1244   |

Table 4: Prompt template for different tasks with 2 demonstrations. For Math problems, we also apply the chat template since we use the instruct models (done by applying the function “apply\_chat\_template” on the instruct models’ tokenizer).

| Name   | Template  |
|--|---|
| Sentiment Classification (SST-2, IMDB, Amazon) | <p>Question: {input_1}<br/>Answer: {output_1}</p> <p>Question: {input_2}<br/>Answer: {output_2}</p> <p>Question: {input_query}<br/>Answer:</p>  |
| Commonsense Reasoning (ARC-Easy, CsQA)         | <p>Question: {input_1}<br/>Answer: {output_1}</p> <p>Question: {input_2}<br/>Answer: {output_2}</p> <p>Question: {input_query}<br/>Answer:</p>  |
| Reading Comprehension (SQuAD, SCIQ)            | <p>Support: {support_1}<br/>Question: {input_1}<br/>Answer: {output_1}</p> <p>Support: {support_2}<br/>Question: {input_2}<br/>Answer: {output_2}</p> <p>Support: {support_query}<br/>Question: {input_query}<br/>Answer:</p>   |
| text to SQL (Geo-Query)                        | <p>Question: {input_1}<br/>Answer: {output_1}</p> <p>Question: {input_2}<br/>Answer: {output_2}</p> <p>Question: {input_query}<br/>Answer:</p>  |
| Math (GSM8K, PRM800K)                          | <p>Question: {input_1}<br/>Answer: {output_1}</p> <p>Question: {input_2}<br/>Answer: {output_2}</p> <p>Let’s think step by step. You need to solve the final<br/> ↔ question and answer in the format: \n#### \{result\<br/> Question: {input_query}<br/> Answer:</p> |

Table 5: **Detailed dataset size before and after sampling.** We show the original and sampled size of demonstration set and test set for all dataset we considered.

| Dataset size      | Classification |         |       | Multiple-choice |      |         | Math  |               | Code     | Reading |       |
|-------------------|----------------|---------|-------|-----------------|------|---------|-------|---------------|----------|---------|-------|
|                   | SST-2          | Amazon  | Imdb  | ARC-Easy        | CsQA | PRM800K | GSM8K | GSM-Plus-Mini | GeoQuery | SQuAD   | SCIQ  |
| Sampled demo set  | 1000           | 1000    | 1000  | 1000            | 1000 | 12000   | 7473  | 7473          | 600      | 10000   | 1000  |
| Sampled test set  | 1000           | 1000    | 1000  | 1000            | 1000 | 500     | 1319  | 2400          | 280      | 1000    | 1000  |
| Original demo set | 67300          | 3600000 | 25000 | 2250            | 9740 | 12000   | 7473  | 7473          | 600      | 87600   | 11700 |
| Original test set | 1820           | 400000  | 25000 | 2380            | 1140 | 500     | 1319  | 2400          | 280      | 10600   | 1000  |

Table 6: Performance of 0-shot and 1-shot Baseline in Code and Reading Tasks. When  $k = 1$ , there is only one possible permutation, so we report a single result for both TopK and TopK-Div methods. For Rand and Div approaches, we report the averaged results across ten random seeds. Embedding = all-roberta-large-v1.

| Model           | Dataset         | $K = 0$ | $K = 1$ |       |       |          | $K = 4$ |       |       |          |
|-----------------|-----------------|---------|---------|-------|-------|----------|---------|-------|-------|----------|
|                 |                 | -       | Rand    | Topk  | Div   | Topk-Div | Rand    | Topk  | Div   | Topk-Div |
| Llama-3.1-8B    | Code (Geoquery) | -       | 2.61    | 37.14 | 16.93 | 37.14    | 12.57   | 63.04 | 33.71 | 71.07    |
|                 | Reading (SQuAD) | 42.30   | 68.64   | 67.00 | 67.87 | 67.00    | 75.95   | 73.51 | 75.66 | 73.28    |
| Gemma-2-9B      | Code (Geoquery) | -       | 3.07    | 41.43 | 16.71 | 41.43    | 13.89   | 61.14 | 36.29 | 70.43    |
|                 | Reading (SQuAD) | 37.90   | 71.34   | 69.00 | 70.69 | 69.00    | 77.19   | 74.82 | 77.06 | 75.05    |
| Mistral-7B-v0.3 | Code (Geoquery) | -       | 2.75    | 40.71 | 18.39 | 40.71    | 12.14   | 60.14 | 34.89 | 71.46    |
|                 | Reading (SQuAD) | 30.50   | 69.12   | 66.30 | 67.80 | 66.30    | 76.70   | 75.04 | 75.96 | 74.43    |

For math problem (generation), given the prompt, we generate the answer using greedy decoding. We do not stop the generation process unless the instruct model generates the stop sign itself. We first try to extract the math expression from the following format “##### {expression}”. If failed, we try to extract from the following format “\boxed{expression}”. If both failed, we extract the final math expression from the answer. The report exact match metric.

For each task, the selected examples in TopK and TopK-Div are fixed, and these two methods are tested once. For Rand and Div, where example selection involves randomness, we test with ten random seeds and report the average results.

## C Additional experiments

In this section, we present some addition (supplementary) experiment results for Section 3. This section is structured as follows:

- Appendix C.1 shows the results of different tasks under 0-shot or 1-shot, to justify the effectiveness of in-context examples;
- Appendix C.2 discusses the best subset\_size in Div;
- Appendix C.3 illustrates the gap between different levels of diversity in TopK-Div;
- Appendix C.4 includes more results and discussions for the OOD setting;
- Appendix C.5 contains the detailed experiments that imply the effect of diversity that is beyond coverage.

### C.1 Results of 0/1-shot

To verify whether the model inherently possesses the ability to solve certain tasks, we tested its 0-shot and 1-shot performance on the SQuAD and GeoQuery datasets. For the Reading task, accuracy is calculated only when the output exactly matches the answer, imposing strict format requirements. Consequently, on SQuAD, once the model understood the output format in the 1-shot setting, the absolute performance gap compared to the 4-shot setting was less than 8%. However, on GeoQuery, even after the model grasped the output format via the 1-shot example, the absolute performance gap compared to the 4-shot setting was still over 20%.

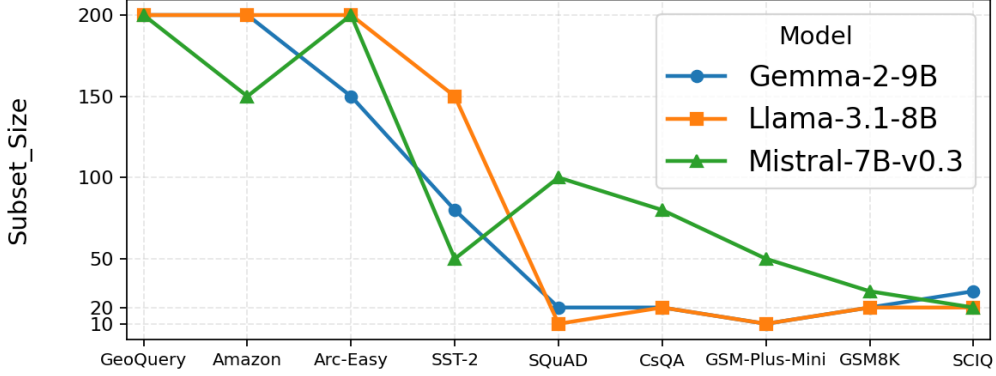


Figure 4: **(Optimal Subset Size of the Div Method for Different Tasks)** We report the optimal subset\_size of the Div method across different tasks. The results show that relatively easier tasks, such as Classification and Multiple-choice, tend to favor larger subset\_size values (with data points concentrated on the left side of the x-axis), whereas more challenging tasks, such as Math and Reading, exhibit substantially smaller optimal subset\_size values, with an overall average not exceeding 30 (with data points concentrated on the right side of the x-axis).

Therefore, the model possesses a strong inherent ability to solve the Reading task (a similar conclusion also holds for the Math task). Conversely, the model itself lacks domain-specific knowledge related to GeoQuery and thus needs to learn more from the provided context.

## C.2 Additional supplement of Finding 1: Adjusting subset sizes for Div

For each task, we identify the subset\_size that yields the best performance for Div ( $subset\_size \in \{10, 20, 30, 50, 80, 100, 120, 150, 200\}$ ). On simpler tasks (e.g., Classification and Multiple-choice), the average optimal subset\_size exceeds 100 (Figure 4), indicating that introducing less diversity is more beneficial. In contrast, on more complex tasks (e.g., Math and Reading), the average optimal subset\_size is below 30, suggesting that incorporating more diversity is advantageous. GeoQuery serves as a typical example of a task that requires strong coverage, for which Div with small subset\_size fails to achieve satisfactory performance.

## C.3 Additional supplement of Finding 1: Adjusting $\alpha$ for TopK-Div

Let  $Acc_\alpha$  denotes the accuracy of TopK-Div parameterized by  $\alpha$  (Equation (3)), We define:

$$\Delta = \frac{1}{5} \sum_{i=6}^{10} Acc_{i/10} - \frac{1}{5} \sum_{i=1}^5 Acc_{i/10}. \quad (4)$$

The difference  $\Delta$  quantifies the gap between the average accuracy of lower diversity and higher diversity in TopK-Div (e.g,  $\Delta > 0$  means less diversity is better). As shown in Table 7, minimal diversity is optimal for simpler tasks, while higher diversity consistently enhances performance as task difficulty increases.

## C.4 More results on OOD setting

In this part, we show the OOD results of math problems on Llama-3.1-8B/70B and Gemma-2-9B/27B instruct-tuned models. We use GSM8K as the demonstration set and PRM800K (Lightman et al., 2023) as the query set. Table 8 summarizes our result. We observe that diversity-aware methods are more robust to this distribution shift. Even for Gemma models where TopK performs very well on ID tasks (demonstration and query set are all PRM800K), TopK is outperformed by diversity-aware methods on the OOD setting. A similar trend also holds for Llama models. One interesting finding is that for PRM800K, more demonstration might not lead to better performance, and also in our experiment, using GSM8K as demonstration works better than using PRM800K data as demonstrations.

Table 7: **Comparison of TopK-Div results with different  $\alpha$ .** We report  $\Delta$  (Equation (4)) for Classification, Multiple-choice and Reading task across six datasets. For each value of  $\alpha$  in TopK-Div, we tested ten permutations and calculated the mean. For relatively simple tasks (SST-2, Amazon, ARC-Easy and CsQA), the average value 0.27% of  $\Delta$  indicates that incorporating less diversity is more beneficial. In contrast, for relatively complex tasks (SCIQ, SQuAD), the average value -0.44% of  $\Delta$  suggests that incorporating more diversity is advantageous. Specifically, the average  $\Delta$  for SST-2 is 0.42%, for ARC-Easy is 0.34%, for CsQA is 0.07%, and for SCIQ is -0.53%. This trend is consistent with our understanding of task difficulty.

| Model        | $K$ | $\Delta$ |        |          |        |        |        |
|--------------|-----|----------|--------|----------|--------|--------|--------|
|              |     | SST-2    | Amazon | ARC-Easy | CsQA   | SCIQ   | SQuAD  |
| Llama-3.2-3B | 4   | 0.11%    | 0.12%  | 0.45%    | 0.41%  | -0.80% | 0.40%  |
|              | 8   | 1.05%    | -0.01% | 0.20%    | 0.06%  | -0.74% | -1.01% |
| Gemma-2-2B   | 4   | 0.34%    | 0.50%  | 0.87%    | -0.37% | -0.07% | 0.04%  |
|              | 8   | 0.19%    | 0.27%  | -0.15%   | 0.18%  | -0.49% | -0.82% |

Table 8: **(Comparison of different methods on math when demonstration and query come from different distribution)** OOD setting for math problem where the test dataset is chosen to be PRM800K. We find that, diversity-aware methods are more superior than TopK in OOD setting. The method that achieves the best in each setting is highlighted.

| Model         | Shots   | Demo.   | Rand         | TopK         | Div          | TopK-Div     |
|---------------|---------|---------|--------------|--------------|--------------|--------------|
| Llama-3.1-8B  | $K = 4$ | PRM800K | 43.50        | 41.40        | <b>44.86</b> | 44.80        |
|               |         | GSM8K   | 41.50        | 41.00        | <b>43.28</b> | 42.00        |
|               | $K = 8$ | PRM800K | 43.32        | 43.00        | <b>44.28</b> | 40.00        |
|               |         | GSM8K   | 41.66        | 42.00        | <b>43.46</b> | 40.80        |
| Llama-3.1-70B | $K = 4$ | PRM800K | 57.78        | 58.20        | 57.42        | <b>59.40</b> |
|               |         | GSM8K   | 60.62        | <b>62.00</b> | 61.88        | 61.00        |
|               | $K = 8$ | PRM800K | 54.72        | <b>59.00</b> | 55.86        | 58.00        |
|               |         | GSM8K   | <b>61.14</b> | 59.60        | 60.96        | 60.00        |
| Gemma-2-9B    | $K = 4$ | PRM800K | 38.04        | 42.40        | 36.78        | <b>44.40</b> |
|               |         | GSM8K   | <b>42.10</b> | 41.00        | 41.04        | 42.20        |
|               | $K = 8$ | PRM800K | 40.66        | <b>46.20</b> | 39.20        | 44.40        |
|               |         | GSM8K   | 42.06        | 41.80        | <b>42.74</b> | 41.60        |
| Gemma-2-27B   | $K = 4$ | PRM800K | 46.06        | 49.20        | 47.80        | <b>49.60</b> |
|               |         | GSM8K   | 46.06        | 46.00        | <b>46.30</b> | 45.20        |
|               | $K = 8$ | PRM800K | 47.10        | <b>50.40</b> | 47.04        | 48.60        |
|               |         | GSM8K   | 45.40        | 44.80        | <b>45.92</b> | 45.20        |

## C.5 Results on perturbation of datasets

To explore whether the way diversity works is by achieving better coverage, we noticed that even when  $k = 1$ , TopK still underperforms Rand/Div methods. We speculate this is because the support in the original dataset contains a lot of noise, causing similar examples not only to fail to provide effective information but also potentially to mislead the model into focusing on noisy information (“coverage” isn’t helpful in such case).

Using DeepSeek-R1, we removed information irrelevant to the answer from the support passages in SQuAD, reducing content by approximately 50%. Based on this, we constructed two variants: SQuAD-Cut, where only the training set is streamlined, and SQuAD-Both-Cut, where both the training and test sets are streamlined. As shown in Table 9, the more streamlined (i.e., higher-quality and less noisy) the dataset, the better the performance of TopK and TopK-Div. Notably, their improvement margins are significantly larger than that of Div (though still more than 1% lower than Div). This indicates

Table 9: Results for SQuAD with cut perturbation. We performed content trimming on the support portion of the SQuAD dataset using Deepseek-r1, retaining only the top 1/3 most answer-relevant content. SQuAD-Cut refers to trimming applied solely to the testing set, while SQuAD-Both-Cut indicates trimming performed on both testing and training sets. The values in parentheses represent performance improvements relative to the original SQuAD dataset.

| Model           | $K$ | Dataset           | Method               |                      |                      |                      |
|-----------------|-----|-------------------|----------------------|----------------------|----------------------|----------------------|
|                 |     |                   | Rand                 | Topk                 | Div                  | Topk-Div             |
| Llama-3.1-8B    | 1   | SQuAD             | 68.64                | 67.00                | 67.87                | 67.00                |
|                 |     | SQuAD-Cut         | 69.43 (+0.79)        | <b>68.20 (+1.20)</b> | 68.45 (+0.58)        | 67.70 (+0.70)        |
|                 |     | SQuAD-Both-Cut    | 69.71 (+1.07)        | <b>69.90 (+2.90)</b> | 69.47 (+1.60)        | <b>69.90 (+2.90)</b> |
|                 | 4   | SQuAD             | 75.95                | 73.51                | 75.66                | 73.28                |
|                 |     | SQuAD-Cut         | 77.15 (+1.2)         | 75.96 (+2.45)        | 77.00 (+1.34)        | <b>76.89 (+2.61)</b> |
|                 |     | SQuAD-Both-Cut    | 76.95 (+1.00)        | 76.15 (+2.64)        | 77.76 (+2.10)        | <b>76.47 (+3.19)</b> |
|                 | 8   | SQuAD             | 77.13                | 75.52                | 77.71                | 76.13                |
|                 |     | SQuAD-Cut         | 79.10 (+1.97)        | 77.43 (+1.91)        | <b>79.43 (+1.72)</b> | 78.64 (+2.51)        |
|                 |     | SQuAD-Both-Cut    | 79.39 (+2.26)        | <b>78.66 (+3.14)</b> | 79.26 (+1.55)        | 79.13 (+3.00)        |
| Gemma-2-9B      | 1   | SQuAD             | 71.34                | 69.00                | 70.69                | 69.00                |
|                 |     | SQuAD-Cut         | 72.96 (+1.62)        | <b>71.20 (+2.20)</b> | 72.25 (+1.56)        | 71.10 (+2.10)        |
|                 |     | SQuAD-Both-Cut    | 73.14 (+1.80)        | <b>72.30 (+3.30)</b> | 72.67 (+1.98)        | <b>72.30 (+3.30)</b> |
|                 | 4   | SQuAD             | 77.19                | 74.82                | 77.06                | 75.05                |
|                 |     | SQuAD-Cut         | 78.75 (+1.56)        | 77.64 (+2.82)        | 78.23 (+1.17)        | <b>78.65 (+3.60)</b> |
|                 |     | SQuAD-Both-Cut    | 78.72 (+1.53)        | <b>77.47 (+2.65)</b> | 78.54 (+1.48)        | 76.78 (+1.73)        |
|                 | 8   | SQuAD             | 79.23                | 77.59                | 79.05                | 77.64                |
|                 |     | SQuAD-Cut         | 80.41 (+1.18)        | 79.74 (+2.15)        | 80.45 (+1.40)        | <b>80.72 (+3.08)</b> |
|                 |     | SQuAD-Both-Cut    | 80.22 (+0.99)        | <b>79.05 (+1.46)</b> | 80.10 (+1.05)        | 78.84 (+1.20)        |
| Mistral-7B-v0.3 | 1   | SQuAD             | 69.12                | 66.30                | 67.80                | 66.30                |
|                 |     | SQuAD-Cut         | 71.38 (+2.26)        | <b>69.70 (+3.40)</b> | 69.21 (+1.41)        | <b>69.70 (+3.40)</b> |
|                 |     | SQuAD-Both-Cut    | 71.44 (+2.32)        | <b>70.70 (+4.40)</b> | 72.00 (+4.20)        | <b>70.70 (+4.40)</b> |
|                 | 4   | SQuAD             | 76.70                | 75.04                | 75.96                | 74.43                |
|                 |     | SQuAD-Cut         | 77.78 (+1.08)        | 77.78 (+2.74)        | 77.18 (+1.22)        | <b>78.70 (+4.37)</b> |
|                 |     | SQuAD-Both-Cut    | 77.76 (+1.06)        | 77.92 (+2.88)        | 77.89 (+1.93)        | <b>77.40 (+2.97)</b> |
|                 | 8   | SQuAD             | 77.30                | 77.05                | 77.67                | 77.44                |
|                 |     | SQuAD-Cut-R1      | 79.00 (+1.70)        | <b>79.09 (+2.04)</b> | 78.64 (+0.97)        | 79.07 (+1.63)        |
|                 |     | SQuAD-Both-Cut-R1 | <b>79.92 (+2.62)</b> | 78.76 (+1.71)        | 79.61 (+1.94)        | 79.64 (+2.20)        |

that when the dataset quality is higher, the ‘‘Coverage’’ mechanism can focus on high signal-to-noise ratio information (rather than incorrectly covering noise), and its effectiveness is significantly enhanced. TopK-based methods are more likely to ‘‘cover’’ high-quality information segments truly relevant to the answer, whereas Div, as an intrinsic metric, inherently includes effective mechanisms not directly dependent on precise semantic coverage (e.g., structural diversity: selecting examples with different sentence structures or argumentation styles). These mechanisms already play a role in the original noisy data, avoiding overfitting to noise, causing it to outperform noise-sensitive coverage strategies, and its baseline performance is already relatively robust. This fully demonstrates that the value of *diversity* is ‘‘beyond coverage’’.

## D Additional ablation studies

### D.1 Results on more models

We evaluated different model sizes from the Gemma and Llama families, including Llama-3.2-1B, Gemma-2-2B, Llama-3.2-3B, Llama-3.1-8B, Gemma-2-9B, Gemma-2-27B, and Llama-3.1-70B. For math tasks, we used the instruct version of the corresponding models. For other tasks, we used the base models. For code tasks, we also tested domain-specific CodeLlama models, including CodeLlama-7B-hf, CodeLlama-13B-hf, and CodeLlama-34B-hf. The results on CodeLlama were consistent with those of other base models.

We report the complete experimental results of the Llama family in Table 11, the Gemma family results in Table 12, and the CodeLlama results in Table 13. The methods that performed well on the corresponding tasks in Table 1 also demonstrated good performance across different model sizes.

Table 10: We supplemented the content omitted in Table 1. The main numerical values represent the mean results over ten random seeds, while the subscript indicates their std. We still highlight the result with the highest mean in bold. In most cases, the fluctuations within each method do not affect our conclusions.

| Model           | $K$        | Method   | Classification               |                              | Multiple-choice              |                              | GSM8K                        | Math<br>GSM-Plus-Mini        | Code<br>GeoQuery             | Reading                      |                              |
|-----------------|------------|----------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|                 |            |          | SST-2                        | Amazon                       | ARC-Easy                     | CsQA                         |                              |                              |                              | SQuAD                        | SCIQ                         |
| Llama-3.1-8B    | 0          | -        | 87.50                        | 95.40                        | 82.43                        | 62.80                        | 53.45                        | 65.12                        | —                            | 42.30                        | 36.40                        |
|                 | 4          | Rand     | 91.31 <sub>0.59</sub>        | <b>96.38</b> <sub>0.25</sub> | 84.72 <sub>0.35</sub>        | 71.15 <sub>0.65</sub>        | <b>82.24</b> <sub>0.52</sub> | 66.90 <sub>0.59</sub>        | 12.57 <sub>1.33</sub>        | <b>75.95</b> <sub>0.55</sub> | 74.00 <sub>0.57</sub>        |
|                 |            | TopK     | <b>94.13</b> <sub>0.21</sub> | 96.24 <sub>0.16</sub>        | <b>86.10</b> <sub>0.32</sub> | 72.54 <sub>0.36</sub>        | 81.99 <sub>0.55</sub>        | 65.30 <sub>0.46</sub>        | 63.04 <sub>1.96</sub>        | 73.51 <sub>0.48</sub>        | 72.70 <sub>0.40</sub>        |
|                 |            | Div      | 91.50 <sub>0.63</sub>        | 96.18 <sub>0.25</sub>        | 85.06 <sub>0.27</sub>        | 71.17 <sub>0.42</sub>        | 82.14 <sub>0.45</sub>        | <b>66.92</b> <sub>0.52</sub> | 33.71 <sub>1.35</sub>        | 75.66 <sub>0.97</sub>        | <b>74.47</b> <sub>0.62</sub> |
|                 | 8          | TopK-Div | 92.75 <sub>0.33</sub>        | 96.15 <sub>0.22</sub>        | 85.83 <sub>0.38</sub>        | <b>72.57</b> <sub>0.35</sub> | 81.74 <sub>0.53</sub>        | 66.12 <sub>0.85</sub>        | <b>71.07</b> <sub>1.11</sub> | 73.28 <sub>0.79</sub>        | 73.87 <sub>0.35</sub>        |
|                 |            | Rand     | 92.27 <sub>0.55</sub>        | <b>96.63</b> <sub>0.27</sub> | 84.38 <sub>0.34</sub>        | 72.23 <sub>0.34</sub>        | 82.81 <sub>0.61</sub>        | <b>66.72</b> <sub>0.72</sub> | 23.21 <sub>1.41</sub>        | 77.13 <sub>0.80</sub>        | 74.65 <sub>0.88</sub>        |
|                 |            | TopK     | <b>93.64</b> <sub>0.36</sub> | 96.12 <sub>0.09</sub>        | <b>85.91</b> <sub>0.29</sub> | <b>73.91</b> <sub>0.38</sub> | 82.26 <sub>0.65</sub>        | 65.99 <sub>0.60</sub>        | 72.04 <sub>0.93</sub>        | 75.52 <sub>0.43</sub>        | 74.72 <sub>0.65</sub>        |
|                 | 8          | Div      | 92.95 <sub>0.35</sub>        | 96.25 <sub>0.19</sub>        | 84.97 <sub>0.32</sub>        | 72.77 <sub>0.61</sub>        | <b>82.98</b> <sub>0.34</sub> | 66.56 <sub>0.60</sub>        | 38.54 <sub>0.90</sub>        | <b>77.71</b> <sub>0.80</sub> | <b>75.17</b> <sub>0.53</sub> |
|                 |            | TopK-Div | 93.33 <sub>0.36</sub>        | 96.43 <sub>0.09</sub>        | 85.39 <sub>0.40</sub>        | 73.76 <sub>0.37</sub>        | 82.63 <sub>0.57</sub>        | 66.48 <sub>0.52</sub>        | <b>78.36</b> <sub>1.24</sub> | 76.13 <sub>0.42</sub>        | 75.07 <sub>0.49</sub>        |
|                 | Gemma-2-9B | 0        | -                            | 67.50                        | 85.10                        | 88.15                        | 61.80                        | 16.07                        | 32.79                        | —                            | 37.90                        |
| 4               |            | Rand     | 93.33 <sub>0.52</sub>        | 96.15 <sub>0.23</sub>        | 89.52 <sub>0.25</sub>        | 74.70 <sub>0.70</sub>        | 84.29 <sub>0.43</sub>        | 74.40 <sub>0.47</sub>        | 13.89 <sub>1.67</sub>        | <b>77.19</b> <sub>0.89</sub> | 75.80 <sub>0.54</sub>        |
|                 |            | TopK     | <b>94.13</b> <sub>0.48</sub> | 96.34 <sub>0.20</sub>        | <b>90.50</b> <sub>0.16</sub> | 75.19 <sub>0.25</sub>        | 84.25 <sub>0.73</sub>        | <b>74.50</b> <sub>0.55</sub> | 61.14 <sub>1.33</sub>        | 74.82 <sub>0.70</sub>        | 75.24 <sub>0.34</sub>        |
|                 |            | Div      | 93.45 <sub>0.46</sub>        | 95.69 <sub>0.23</sub>        | 90.03 <sub>0.24</sub>        | 74.85 <sub>0.39</sub>        | <b>84.44</b> <sub>0.91</sub> | 73.34 <sub>0.62</sub>        | 36.29 <sub>1.05</sub>        | 77.06 <sub>0.57</sub>        | <b>75.96</b> <sub>0.55</sub> |
| 8               |            | TopK-Div | 93.34 <sub>0.34</sub>        | <b>96.57</b> <sub>0.16</sub> | 90.19 <sub>0.19</sub>        | <b>75.60</b> <sub>0.54</sub> | 83.54 <sub>0.56</sub>        | 74.47 <sub>0.63</sub>        | <b>70.43</b> <sub>1.24</sub> | 75.05 <sub>0.41</sub>        | 75.21 <sub>0.29</sub>        |
|                 |            | Rand     | 93.30 <sub>0.36</sub>        | 96.09 <sub>0.23</sub>        | 89.39 <sub>0.28</sub>        | 75.98 <sub>0.56</sub>        | <b>84.34</b> <sub>0.54</sub> | 74.48 <sub>0.63</sub>        | 24.36 <sub>1.19</sub>        | <b>79.23</b> <sub>0.64</sub> | 76.28 <sub>0.50</sub>        |
|                 |            | TopK     | <b>94.20</b> <sub>0.28</sub> | 96.55 <sub>0.16</sub>        | <b>90.62</b> <sub>0.16</sub> | 76.14 <sub>0.63</sub>        | 83.57 <sub>0.53</sub>        | <b>75.36</b> <sub>0.43</sub> | 71.00 <sub>1.20</sub>        | 77.59 <sub>0.42</sub>        | 75.55 <sub>0.18</sub>        |
| 8               |            | Div      | 93.41 <sub>0.20</sub>        | 95.94 <sub>0.25</sub>        | 89.90 <sub>0.19</sub>        | <b>76.60</b> <sub>0.32</sub> | 84.22 <sub>0.52</sub>        | 74.69 <sub>0.64</sub>        | 42.07 <sub>1.10</sub>        | 79.05 <sub>0.93</sub>        | <b>76.65</b> <sub>0.60</sub> |
|                 |            | TopK-Div | 94.04 <sub>0.29</sub>        | <b>96.58</b> <sub>0.04</sub> | 90.48 <sub>0.22</sub>        | 76.53 <sub>0.21</sub>        | 83.85 <sub>0.66</sub>        | 75.16 <sub>0.32</sub>        | <b>76.32</b> <sub>0.85</sub> | 77.64 <sub>0.63</sub>        | 76.24 <sub>0.48</sub>        |
| Mistral-7B-v0.3 |            | 0        | -                            | 66.50                        | 94.00                        | 76.41                        | 51.80                        | 9.48                         | 5.17                         | —                            | 30.50                        |
|                 | 4          | Rand     | 91.00 <sub>0.78</sub>        | 94.02 <sub>0.61</sub>        | 82.77 <sub>0.48</sub>        | 69.83 <sub>0.81</sub>        | 48.78 <sub>1.00</sub>        | 37.20 <sub>0.69</sub>        | 12.14 <sub>1.47</sub>        | <b>76.70</b> <sub>0.72</sub> | 74.71 <sub>0.54</sub>        |
|                 |            | TopK     | <b>93.57</b> <sub>0.25</sub> | <b>96.17</b> <sub>0.20</sub> | <b>85.21</b> <sub>0.30</sub> | 69.73 <sub>0.43</sub>        | 49.28 <sub>1.17</sub>        | 38.20 <sub>0.55</sub>        | 60.14 <sub>0.82</sub>        | 75.04 <sub>0.74</sub>        | 73.73 <sub>0.59</sub>        |
|                 |            | Div      | 91.98 <sub>0.46</sub>        | 94.15 <sub>0.31</sub>        | 82.98 <sub>0.25</sub>        | <b>70.15</b> <sub>0.56</sub> | 49.49 <sub>0.87</sub>        | 37.50 <sub>0.76</sub>        | 34.89 <sub>1.39</sub>        | 75.96 <sub>1.08</sub>        | <b>75.83</b> <sub>0.57</sub> |
|                 | 8          | TopK-Div | 92.73 <sub>0.30</sub>        | 95.90 <sub>0.15</sub>        | 84.55 <sub>0.20</sub>        | 69.91 <sub>0.49</sub>        | <b>49.99</b> <sub>1.02</sub> | <b>38.45</b> <sub>0.81</sub> | <b>71.46</b> <sub>1.35</sub> | 74.43 <sub>0.50</sub>        | 73.16 <sub>0.28</sub>        |
|                 |            | Rand     | 92.49 <sub>0.34</sub>        | 95.35 <sub>0.36</sub>        | 83.69 <sub>0.36</sub>        | 71.65 <sub>0.60</sub>        | 47.86 <sub>1.19</sub>        | 36.32 <sub>0.71</sub>        | 22.18 <sub>1.96</sub>        | 77.30 <sub>0.54</sub>        | 75.54 <sub>0.63</sub>        |
|                 |            | TopK     | <b>93.61</b> <sub>0.32</sub> | <b>96.15</b> <sub>0.16</sub> | <b>85.17</b> <sub>0.25</sub> | 71.88 <sub>0.38</sub>        | 48.43 <sub>1.02</sub>        | 37.35 <sub>0.53</sub>        | 70.50 <sub>1.36</sub>        | 77.05 <sub>0.41</sub>        | 75.44 <sub>0.43</sub>        |
|                 | 8          | Div      | 92.55 <sub>0.29</sub>        | 95.10 <sub>0.37</sub>        | 84.27 <sub>0.41</sub>        | <b>72.04</b> <sub>0.61</sub> | 48.33 <sub>1.10</sub>        | 36.12 <sub>0.34</sub>        | 39.14 <sub>1.40</sub>        | <b>77.67</b> <sub>1.56</sub> | <b>76.30</b> <sub>0.31</sub> |
|                 |            | TopK-Div | 93.47 <sub>0.41</sub>        | 96.11 <sub>0.16</sub>        | 84.85 <sub>0.34</sub>        | 71.81 <sub>0.19</sub>        | <b>48.60</b> <sub>0.71</sub> | <b>37.81</b> <sub>0.76</sub> | <b>77.93</b> <sub>1.70</sub> | 77.44 <sub>0.37</sub>        | 75.22 <sub>0.42</sub>        |

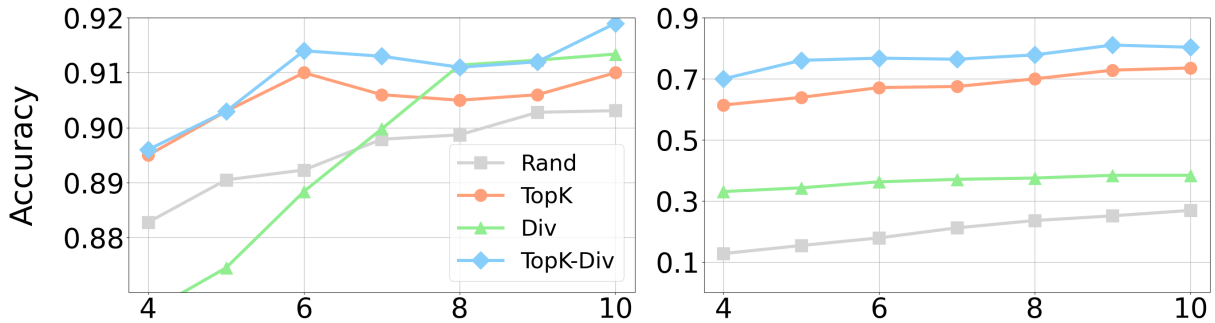


Figure 5: The performance of different demonstration selection methods with different number of shots  $K$ . **Left:** sentiment classification task with demonstrations come from Amazon and queries come from SST-2. **Right:** text to SQL task with demonstrations and query come from the training and test set of GeoQuery Standard Split.

Due to resource constraints, our experiments primarily focused on mainstream open-source models. We tested the Math task on the commercial-grade models gpt-4o-mini and deepseek-v3. As shown in Table 14, the Div method consistently outperformed TopK on both gsm8k and prm800k.

## D.2 Changing the number of shots

In this section, we investigate how the performance advantage of diversity-aware methods over TopK evolves with increasing shot count. Our results in Figure 5 show that the improvement from diversity-aware selection (Div) remains substantial even with higher number of shots.

We believe that as the shot number increases, there is an increase in redundant information among the examples selected by the TopK method. In contrast, the TopK-Div method minimizes the occurrence of redundant information as much as possible, thereby enabling the model to more clearly identify the task theme.

Figure 5 presents the relative improvement on the GeoQuery standard split and SCIQ — two tasks where diversity-aware methods showed clear benefits (Section 3.1) — across different sizes of the Llama-3.1/3.2 and Gemma-2 model families. The results indicate that, in general, the relative improvement from diversity-aware selection does not diminish significantly as model size increases. This underscores the continued importance of understanding diversity’s role in demonstration selection.

Table 11: Performance of different algorithms for models belong to Llama-family. Setting same as Table 1 while adding results from more models (Llama-3.2-1B, Llama-3.2-3B, Llama-3.1-70B). Our finding that diversity helps for more challenging tasks still holds.

| model         | K | Method   | Classification               |                              | Multiple-choice              |                              |                              | Math                         | Code                         | Reading                      |                              |
|---------------|---|----------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|               |   |          | SST-2                        | Amazon                       | Arc-easy                     | CsQA                         | GSM8K                        | GSM-Plus-Mini                | GeoQuery                     | SQuAD                        | SCIQ                         |
| Llama-3.2-1B  | 4 | Rand     | 86.88 <sub>0.49</sub>        | 90.64 <sub>0.43</sub>        | 71.65 <sub>0.46</sub>        | 59.54 <sub>0.63</sub>        | <b>22.18</b> <sub>0.81</sub> | 12.62 <sub>0.56</sub>        | 7.43 <sub>1.27</sub>         | 56.17 <sub>0.75</sub>        | 62.75 <sub>0.84</sub>        |
|               |   | TopK     | <b>91.87</b> <sub>0.49</sub> | <b>93.46</b> <sub>0.25</sub> | <b>75.28</b> <sub>0.57</sub> | 60.20 <sub>0.46</sub>        | 19.57 <sub>0.71</sub>        | 9.41 <sub>0.60</sub>         | 48.25 <sub>0.94</sub>        | 55.09 <sub>0.38</sub>        | <b>63.65</b> <sub>0.36</sub> |
|               |   | Div      | 88.35 <sub>1.19</sub>        | 91.14 <sub>0.43</sub>        | 73.52 <sub>0.54</sub>        | 60.60 <sub>0.65</sub>        | 21.29 <sub>1.05</sub>        | <b>13.16</b> <sub>0.73</sub> | 29.46 <sub>1.70</sub>        | 55.40 <sub>1.18</sub>        | 63.38 <sub>0.54</sub>        |
|               | 8 | TopK-Div | 91.47 <sub>0.60</sub>        | 93.22 <sub>0.33</sub>        | 74.62 <sub>0.44</sub>        | <b>60.89</b> <sub>0.36</sub> | 20.30 <sub>0.80</sub>        | 9.35 <sub>0.53</sub>         | <b>56.57</b> <sub>1.18</sub> | <b>56.39</b> <sub>0.96</sub> | 62.96 <sub>0.31</sub>        |
|               |   | Rand     | 89.56 <sub>0.81</sub>        | 92.62 <sub>0.39</sub>        | 72.72 <sub>0.32</sub>        | 61.27 <sub>0.62</sub>        | 21.04 <sub>0.63</sub>        | 10.01 <sub>0.44</sub>        | 13.04 <sub>1.71</sub>        | <b>58.76</b> <sub>0.56</sub> | 65.05 <sub>0.45</sub>        |
|               |   | TopK     | <b>92.91</b> <sub>0.30</sub> | 93.95 <sub>0.24</sub>        | <b>75.41</b> <sub>0.40</sub> | 61.77 <sub>0.35</sub>        | 16.58 <sub>0.53</sub>        | 7.12 <sub>0.44</sub>         | 56.29 <sub>1.94</sub>        | 58.38 <sub>0.66</sub>        | 65.87 <sub>0.45</sub>        |
| Llama-3.2-3B  | 4 | Div      | 87.96 <sub>1.14</sub>        | 92.72 <sub>0.34</sub>        | 73.53 <sub>0.45</sub>        | <b>62.24</b> <sub>0.55</sub> | <b>22.24</b> <sub>0.93</sub> | <b>11.73</b> <sub>1.55</sub> | 32.75 <sub>1.46</sub>        | 58.37 <sub>1.16</sub>        | 65.89 <sub>0.94</sub>        |
|               |   | TopK-Div | 92.30 <sub>0.38</sub>        | <b>94.06</b> <sub>0.20</sub> | 74.74 <sub>0.28</sub>        | 62.20 <sub>0.50</sub>        | 16.00 <sub>0.81</sub>        | 6.45 <sub>0.46</sub>         | <b>65.11</b> <sub>1.63</sub> | 58.27 <sub>0.74</sub>        | <b>66.24</b> <sub>0.57</sub> |
|               |   | Rand     | 90.40 <sub>0.66</sub>        | 95.87 <sub>0.21</sub>        | 78.62 <sub>0.44</sub>        | 68.51 <sub>0.64</sub>        | 69.64 <sub>0.98</sub>        | 50.50 <sub>0.59</sub>        | 9.86 <sub>1.28</sub>         | <b>71.59</b> <sub>0.60</sub> | 72.43 <sub>0.70</sub>        |
| Llama-3.1-8B  | 4 | TopK     | 92.87 <sub>0.31</sub>        | 96.25 <sub>0.24</sub>        | 81.51 <sub>0.34</sub>        | 68.72 <sub>0.37</sub>        | <b>70.05</b> <sub>0.86</sub> | 50.10 <sub>0.53</sub>        | 54.04 <sub>1.68</sub>        | 71.21 <sub>0.48</sub>        | 70.87 <sub>0.46</sub>        |
|               |   | Div      | 90.87 <sub>0.59</sub>        | 95.57 <sub>0.18</sub>        | 80.41 <sub>0.60</sub>        | 68.80 <sub>0.57</sub>        | 68.71 <sub>0.75</sub>        | <b>51.53</b> <sub>0.87</sub> | 31.21 <sub>1.82</sub>        | 71.13 <sub>1.42</sub>        | <b>72.58</b> <sub>0.61</sub> |
|               |   | TopK-Div | <b>93.03</b> <sub>0.29</sub> | <b>96.43</b> <sub>0.15</sub> | <b>81.71</b> <sub>0.30</sub> | <b>68.95</b> <sub>0.25</sub> | 69.14 <sub>0.86</sub>        | 50.60 <sub>0.38</sub>        | <b>59.75</b> <sub>2.03</sub> | 70.73 <sub>0.49</sub>        | 71.28 <sub>0.40</sub>        |
|               | 8 | Rand     | 91.79 <sub>0.36</sub>        | 96.09 <sub>0.14</sub>        | 78.91 <sub>0.40</sub>        | 69.89 <sub>0.68</sub>        | 68.79 <sub>0.94</sub>        | <b>51.12</b> <sub>0.81</sub> | 19.29 <sub>1.60</sub>        | 73.14 <sub>0.63</sub>        | 72.57 <sub>0.85</sub>        |
|               |   | TopK     | <b>93.71</b> <sub>0.40</sub> | 96.18 <sub>0.20</sub>        | 81.03 <sub>0.36</sub>        | <b>70.53</b> <sub>0.33</sub> | <b>69.40</b> <sub>0.89</sub> | 50.00 <sub>0.84</sub>        | 61.89 <sub>1.49</sub>        | 72.61 <sub>0.30</sub>        | 71.83 <sub>0.46</sub>        |
|               |   | Div      | 91.99 <sub>0.47</sub>        | 95.83 <sub>0.27</sub>        | 80.68 <sub>0.34</sub>        | 70.26 <sub>0.40</sub>        | 66.54 <sub>1.14</sub>        | 50.87 <sub>0.61</sub>        | 36.71 <sub>1.39</sub>        | 72.76 <sub>1.53</sub>        | <b>74.07</b> <sub>0.49</sub> |
| Llama-3.1-70B | 4 | TopK-Div | 93.55 <sub>0.35</sub>        | <b>96.37</b> <sub>0.17</sub> | <b>81.57</b> <sub>0.30</sub> | 70.18 <sub>0.33</sub>        | 69.38 <sub>0.98</sub>        | 49.96 <sub>0.73</sub>        | <b>72.11</b> <sub>1.77</sub> | <b>73.80</b> <sub>0.56</sub> | 72.08 <sub>0.53</sub>        |
|               |   | Rand     | 91.31 <sub>0.59</sub>        | <b>96.38</b> <sub>0.25</sub> | 84.72 <sub>0.35</sub>        | 71.15 <sub>0.65</sub>        | <b>82.24</b> <sub>0.52</sub> | 66.90 <sub>0.59</sub>        | 12.57 <sub>1.33</sub>        | <b>75.95</b> <sub>0.55</sub> | 74.00 <sub>0.57</sub>        |
|               |   | TopK     | <b>94.13</b> <sub>0.21</sub> | 96.24 <sub>0.16</sub>        | <b>86.10</b> <sub>0.32</sub> | 72.54 <sub>0.36</sub>        | 81.99 <sub>0.55</sub>        | 65.30 <sub>0.46</sub>        | 63.04 <sub>1.96</sub>        | 73.51 <sub>0.48</sub>        | 72.70 <sub>0.40</sub>        |
|               | 8 | Div      | 91.50 <sub>0.63</sub>        | 96.18 <sub>0.25</sub>        | 85.06 <sub>0.27</sub>        | 71.17 <sub>0.42</sub>        | 82.14 <sub>0.45</sub>        | <b>66.92</b> <sub>0.52</sub> | 33.71 <sub>1.35</sub>        | 75.66 <sub>0.97</sub>        | <b>74.47</b> <sub>0.62</sub> |
|               |   | TopK-Div | 92.75 <sub>0.33</sub>        | 96.15 <sub>0.22</sub>        | 85.83 <sub>0.38</sub>        | <b>72.57</b> <sub>0.35</sub> | 81.74 <sub>0.53</sub>        | 66.12 <sub>0.85</sub>        | <b>71.07</b> <sub>1.11</sub> | 73.28 <sub>0.79</sub>        | 73.87 <sub>0.35</sub>        |
|               |   | Rand     | 92.27 <sub>0.55</sub>        | <b>96.63</b> <sub>0.27</sub> | 84.38 <sub>0.34</sub>        | 72.23 <sub>0.34</sub>        | 82.81 <sub>0.61</sub>        | <b>66.72</b> <sub>0.72</sub> | 23.21 <sub>1.41</sub>        | 77.13 <sub>0.80</sub>        | 74.65 <sub>0.88</sub>        |
| Llama-3.1-70B | 4 | TopK     | <b>93.64</b> <sub>0.36</sub> | 96.12 <sub>0.09</sub>        | <b>85.91</b> <sub>0.29</sub> | <b>73.91</b> <sub>0.38</sub> | 82.26 <sub>0.65</sub>        | 65.99 <sub>0.60</sub>        | 72.04 <sub>0.93</sub>        | 75.52 <sub>0.43</sub>        | 74.72 <sub>0.65</sub>        |
|               |   | Div      | 92.95 <sub>0.35</sub>        | 96.25 <sub>0.19</sub>        | 84.97 <sub>0.32</sub>        | 72.77 <sub>0.61</sub>        | <b>82.98</b> <sub>0.34</sub> | 66.56 <sub>0.60</sub>        | 38.54 <sub>0.90</sub>        | <b>77.71</b> <sub>0.80</sub> | <b>75.17</b> <sub>0.53</sub> |
|               |   | TopK-Div | 93.33 <sub>0.36</sub>        | 96.43 <sub>0.09</sub>        | 85.39 <sub>0.40</sub>        | 73.76 <sub>0.37</sub>        | 82.63 <sub>0.57</sub>        | 66.48 <sub>0.52</sub>        | <b>78.36</b> <sub>1.24</sub> | 76.13 <sub>0.42</sub>        | 75.07 <sub>0.49</sub>        |
|               | 8 | Rand     | 94.16 <sub>0.33</sub>        | 96.77 <sub>0.38</sub>        | 89.76 <sub>0.16</sub>        | 75.48 <sub>0.62</sub>        | 88.64 <sub>0.48</sub>        | 77.14 <sub>0.39</sub>        | 17.50 <sub>1.88</sub>        | <b>81.47</b> <sub>0.75</sub> | 75.51 <sub>0.87</sub>        |
|               |   | TopK     | <b>94.81</b> <sub>0.34</sub> | 96.86 <sub>0.11</sub>        | <b>90.57</b> <sub>0.28</sub> | <b>76.22</b> <sub>0.28</sub> | 88.87 <sub>0.52</sub>        | 76.19 <sub>0.57</sub>        | 66.46 <sub>1.23</sub>        | 79.15 <sub>0.22</sub>        | 75.67 <sub>0.38</sub>        |
|               |   | Div      | 94.34 <sub>0.28</sub>        | 96.36 <sub>0.23</sub>        | 90.14 <sub>0.30</sub>        | 75.53 <sub>0.33</sub>        | <b>89.27</b> <sub>0.53</sub> | <b>77.21</b> <sub>0.47</sub> | 39.00 <sub>1.87</sub>        | 81.27 <sub>1.25</sub>        | <b>77.75</b> <sub>0.44</sub> |
| Llama-3.1-70B | 4 | TopK-Div | 94.20 <sub>0.18</sub>        | <b>96.88</b> <sub>0.12</sub> | 90.46 <sub>0.32</sub>        | 76.21 <sub>0.49</sub>        | 88.67 <sub>0.42</sub>        | 76.94 <sub>0.45</sub>        | <b>77.32</b> <sub>0.95</sub> | 79.26 <sub>0.37</sub>        | 75.59 <sub>0.33</sub>        |
|               |   | Rand     | 94.66 <sub>0.36</sub>        | 96.95 <sub>0.28</sub>        | 89.84 <sub>0.28</sub>        | 77.14 <sub>0.40</sub>        | 89.47 <sub>0.54</sub>        | 76.93 <sub>0.77</sub>        | 26.89 <sub>1.90</sub>        | 82.62 <sub>0.47</sub>        | 76.56 <sub>0.78</sub>        |
|               |   | TopK     | 94.18 <sub>0.32</sub>        | 96.95 <sub>0.18</sub>        | 90.24 <sub>0.25</sub>        | <b>77.65</b> <sub>0.30</sub> | 89.33 <sub>0.29</sub>        | 76.29 <sub>0.46</sub>        | 75.68 <sub>1.08</sub>        | 81.38 <sub>0.51</sub>        | 76.70 <sub>0.63</sub>        |
| Llama-3.1-70B | 8 | Div      | <b>94.95</b> <sub>0.30</sub> | 96.47 <sub>0.25</sub>        | 89.99 <sub>0.17</sub>        | 77.33 <sub>0.59</sub>        | <b>89.65</b> <sub>0.27</sub> | <b>77.11</b> <sub>0.29</sub> | 44.25 <sub>1.92</sub>        | <b>83.18</b> <sub>1.43</sub> | <b>78.37</b> <sub>0.60</sub> |
|               |   | TopK-Div | 94.75 <sub>0.19</sub>        | <b>97.27</b> <sub>0.11</sub> | <b>90.71</b> <sub>0.25</sub> | 77.24 <sub>0.31</sub>        | 89.17 <sub>0.42</sub>        | 76.74 <sub>0.92</sub>        | <b>81.39</b> <sub>1.16</sub> | 81.47 <sub>0.23</sub>        | 76.77 <sub>0.41</sub>        |

Table 12: Performance of different algorithms for models belong to Gemma-family. Setting same as Table 1 while adding results from more models (Gemma-2-2b and Gemma-2-27b). Our finding that diversity helps for more challenging tasks still holds.

| model       | K | Method   | Classification               |                              | Multiple-choice              |                              |                              | Math                         | Code                         | Reading                      |                              |
|-------------|---|----------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|             |   |          | SST-2                        | Amazon                       | Arc-easy                     | CsQA                         | GSM8K                        | GSM-Plus-Mini                | GeoQuery                     | SQuAD                        | SCIQ                         |
| Gemma-2-2B  | 4 | Rand     | 85.00 <sub>0.94</sub>        | 92.34 <sub>0.75</sub>        | 82.84 <sub>0.45</sub>        | 68.89 <sub>0.61</sub>        | 40.45 <sub>1.11</sub>        | 33.43 <sub>0.69</sub>        | 9.14 <sub>1.38</sub>         | <b>69.03</b> <sub>0.34</sub> | 71.27 <sub>0.63</sub>        |
|             |   | TopK     | 90.67 <sub>0.55</sub>        | <b>95.14</b> <sub>0.18</sub> | <b>84.57</b> <sub>0.35</sub> | 69.93 <sub>0.47</sub>        | 41.33 <sub>0.85</sub>        | <b>34.39</b> <sub>0.71</sub> | 53.39 <sub>1.88</sub>        | 68.19 <sub>0.76</sub>        | 71.09 <sub>0.49</sub>        |
|             |   | Div      | 89.63 <sub>0.54</sub>        | 92.55 <sub>0.44</sub>        | 84.21 <sub>0.37</sub>        | <b>71.03</b> <sub>0.62</sub> | 40.53 <sub>2.15</sub>        | 32.75 <sub>1.64</sub>        | 31.29 <sub>1.47</sub>        | 67.73 <sub>1.18</sub>        | <b>72.34</b> <sub>0.45</sub> |
|             | 8 | TopK-Div | <b>91.66</b> <sub>0.57</sub> | 95.01 <sub>0.22</sub>        | 84.51 <sub>0.32</sub>        | 70.72 <sub>0.52</sub>        | <b>42.74</b> <sub>0.57</sub> | <b>34.39</b> <sub>0.64</sub> | <b>61.04</b> <sub>1.55</sub> | 67.85 <sub>0.60</sub>        | 71.64 <sub>0.30</sub>        |
|             |   | Rand     | 89.96 <sub>0.51</sub>        | 94.10 <sub>0.47</sub>        | 82.62 <sub>0.33</sub>        | 70.26 <sub>0.34</sub>        | 36.44 <sub>0.82</sub>        | 34.55 <sub>0.46</sub>        | 16.68 <sub>2.01</sub>        | 69.99 <sub>0.68</sub>        | 72.12 <sub>0.29</sub>        |
|             |   | TopK     | 92.22 <sub>0.46</sub>        | <b>95.58</b> <sub>0.23</sub> | 84.30 <sub>0.19</sub>        | 71.06 <sub>0.44</sub>        | <b>43.05</b> <sub>0.69</sub> | 36.45 <sub>0.42</sub>        | 61.00 <sub>1.44</sub>        | 69.15 <sub>0.44</sub>        | <b>72.35</b> <sub>0.44</sub> |
| Gemma-2-9B  | 4 | Div      | 91.81 <sub>0.48</sub>        | 94.57 <sub>0.37</sub>        | 84.37 <sub>0.38</sub>        | <b>72.22</b> <sub>0.62</sub> | 38.87 <sub>1.28</sub>        | 34.07 <sub>1.12</sub>        | 35.29 <sub>1.25</sub>        | 69.31 <sub>0.99</sub>        | 72.28 <sub>0.43</sub>        |
|             |   | TopK-Div | <b>92.40</b> <sub>0.28</sub> | 95.53 <sub>0.20</sub>        | <b>84.39</b> <sub>0.24</sub> | 71.99 <sub>0.36</sub>        | 42.93 <sub>0.66</sub>        | <b>36.47</b> <sub>0.47</sub> | <b>68.93</b> <sub>1.35</sub> | <b>70.09</b> <sub>0.54</sub> | 72.30 <sub>0.37</sub>        |
|             |   | Rand     | 93.33 <sub>0.52</sub>        | 96.15 <sub>0.23</sub>        | 89.52 <sub>0.25</sub>        | 74.70 <sub>0.70</sub>        | 84.29 <sub>0.43</sub>        | 74.40 <sub>0.47</sub>        | 13.89 <sub>1.67</sub>        | <b>77.19</b> <sub>0.89</sub> | 75.80 <sub>0.54</sub>        |
|             | 8 | TopK     | <b>94.47</b> <sub>0.48</sub> | 96.34 <sub>0.20</sub>        | <b>90.50</b> <sub>0.16</sub> | 75.19 <sub>0.25</sub>        | 84.25 <sub>0.73</sub>        | <b>74.50</b> <sub>0.55</sub> | 61.14 <sub>1.33</sub>        | 74.82 <sub>0.70</sub>        | 75.24 <sub>0.34</sub>        |
|             |   | Div      | 93.45 <sub>0.46</sub>        | 95.69 <sub>0.23</sub>        | 90.03 <sub>0.24</sub>        | 74.85 <sub>0.39</sub>        | <b>84.44</b> <sub>0.91</sub> | 73.34 <sub>0.62</sub>        | 36.29 <sub>1.05</sub>        | 77.06 <sub>0.57</sub>        | <b>75.96</b> <sub>0.55</sub> |
|             |   | TopK-Div | 93.34 <sub>0.34</sub>        | <b>96.57</b> <sub>0.16</sub> | 90.19 <sub>0.19</sub>        | <b>75.60</b> <sub>0.54</sub> | 83.54 <sub>0.56</sub>        | 74.47 <sub>0.63</sub>        | <b>70.43</b> <sub>1.24</sub> | 75.05 <sub>0.41</sub>        | 75.21 <sub>0.29</sub>        |
| Gemma-2-27B | 4 | Rand     | 93.30 <sub>0.36</sub>        | 96.09 <sub>0.23</sub>        | 89.39 <sub>0.28</sub>        | 75.98 <sub>0.56</sub>        | <b>84.34</b> <sub>0.54</sub> | 74.48 <sub>0.63</sub>        | 24.36 <sub>1.19</sub>        | <b>79.23</b> <sub>0.64</sub> | 76.28 <sub>0.50</sub>        |
|             |   | TopK     | <b>94.20</b> <sub>0.28</sub> | 96.55 <sub>0.16</sub>        | <b>90.62</b> <sub>0.19</sub> | 76.14 <sub>0.63</sub>        | 83.57 <sub>0.53</sub>        | <b>75.36</b> <sub>0.43</sub> | 71.00 <sub>1.20</sub>        | 77.59 <sub>0.42</sub>        | 75.59 <sub>0.18</sub>        |
|             |   | Div      | 93.41 <sub>0.20</sub>        | 95.94 <sub>0.25</sub>        | 89.90 <sub>0.19</sub>        | <b>76.60</b> <sub>0.32</sub> | 84.22 <sub>0.52</sub>        | 74.69 <sub>0.64</sub>        | 42.07 <sub>1.10</sub>        | 79.05 <sub>0.93</sub>        | <b>76.65</b> <sub>0.60</sub> |
|             | 8 | TopK-Div | 94.04 <sub>0.29</sub>        | <b>96.58</b> <sub>0.04</sub> | 90.48 <sub>0.22</sub>        | 76.53 <sub>0.21</sub>        | 83.85 <sub>0.66</sub>        | 75.16 <sub>0.32</sub>        | <b>76.32</b> <sub>0.85</sub> | 77.64 <sub>0.63</sub>        | 76.24 <sub>0.48</sub>        |
|             |   | Rand     | 94.16 <sub>0.40</sub>        | 96.06 <sub>0.26</sub>        | <b>89.99</b> <sub>0.39</sub> | 76.15 <sub>0.41</sub>        | 90.16 <sub>0.33</sub>        | <b>70.76</b> <sub>0.71</sub> | 18.68 <sub>1.83</sub>        | <b>80.54</b> <sub>0.59</sub> | 75.61 <sub>0.78</sub>        |
|             |   | TopK     | <b>95.00</b> <sub>0.33</sub> | 96.47 <sub>0.11</sub>        | 89.64 <sub>0.15</sub>        | 76.47 <sub>0.46</sub>        | 89.73 <sub>0.38</sub>        | 69.27 <sub>0.49</sub>        | 67.75 <sub>1.45</sub>        | 78.43 <sub>0.49</sub>        | 76.35 <sub>0.46</sub>        |
| Gemma-2-27B | 4 | Div      | 94.15 <sub>0.43</sub>        | 95.57 <sub>0.36</sub>        | 89.78 <sub>0.40</sub>        | <b>76.97</b> <sub>0.47</sub> | <b>90.68</b> <sub>0.23</sub> | 69.85 <sub>1.60</sub>        | 41.75 <sub>2.16</sub>        | 79.91 <sub>1.14</sub>        | <b>76.73</b> <sub>0.61</sub> |
|             |   | TopK-Div | 94.43 <sub>0.18</sub>        | <b>96.59</b> <sub>0.12</sub> | 89.76 <sub>0.16</sub>        | 76.15 <sub>0.41</sub>        | 89.53 <sub>0.23</sub>        | 69.62 <sub>0.56</sub>        | <b>79.11</b> <sub>0.97</sub> | 78.39 <sub>0.46</sub>        | 75.91 <sub>0.53</sub>        |
|             |   | Rand     | 94.59 <sub>0.45</sub>        | 96.42 <sub>0.46</sub>        | 89.62 <sub>0.21</sub>        | 77.17 <sub>0.69</sub>        | 90.23 <sub>0.25</sub>        | 69.04 <sub>0.53</sub>        | 30.36 <sub>2.04</sub>        | <b>81.81</b> <sub>0.37</sub> | 77.09 <sub>0.53</sub>        |
|             | 8 | TopK     | 94.30 <sub>0.32</sub>        | <b>96.60</b> <sub>0.18</sub> | 90.14 <sub>0.24</sub>        | 77.20 <sub>0.42</sub>        | 89.95 <sub>0.27</sub>        | 66.54 <sub>0.37</sub>        | 77.54 <sub>1.00</sub>        | 80.72 <sub>0.38</sub>        | 76.42 <sub>0.52</sub>        |
|             |   | Div      | <b>94.61</b> <sub>0.29</sub> | 95.95 <sub>0.27</sub>        | 90.21 <sub>0.23</sub>        | <b>78.47</b> <sub>0.41</sub> | <b>90.45</b> <sub>0.22</sub> | <b>70.02</b> <sub>0.91</sub> | 46.96 <sub>2.23</sub>        | 81.38 <sub>0.89</sub>        | <b>77.84</b> <sub>0.51</sub> |
|             |   | TopK-Div | 94.56 <sub>0.29</sub>        | 96.43 <sub>0.13</sub>        | <b>90.34</b> <sub>0.22</sub> | 77.29 <sub>0.31</sub>        | 89.70 <sub>0.25</sub>        | 68.48 <sub>0.38</sub>        | <b>82.</b>                   |                              |                              |

Table 13: **CodeLlama-family results on GeoQuery dataset with different split.** We observe that on GeoQuery dataset, TopK-Div consistently works better than TopK, and there is also a large gap between TopK and more diversity-aware methods like Div and Rand, which aligns with the results in Table 11, Table 12, and Table 1 for Mistral-v0.3. The gap between different methods is wide and std is small, so we omit the std.

| model            | $K$          | Method       | GeoQuery     |              |              |              |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                  |              |              | Standard     | Tmcd         | Template     | Length       |
| CodeLlama-7B-hf  | 4            | Rand         | 12.21        | 10.43        | 9.75         | 3.61         |
|                  |              | TopK         | 57.86        | 35.68        | 36.90        | 25.91        |
|                  |              | Div          | 33.11        | 21.95        | 27.22        | 13.16        |
|                  |              | TopK-Div     | <b>67.86</b> | <b>40.00</b> | <b>50.34</b> | <b>33.64</b> |
|                  | 8            | Rand         | 21.11        | 17.25        | 17.93        | 8.05         |
|                  |              | TopK         | 58.21        | 42.95        | 48.06        | 32.95        |
| Div              |              | 38.29        | 24.75        | 31.41        | 16.48        |              |
| TopK-Div         | <b>66.79</b> | <b>46.36</b> | <b>55.13</b> | <b>39.09</b> |              |              |
| CodeLlama-13B-hf | 4            | Rand         | 13.82        | 11.66        | 11.73        | 4.11         |
|                  |              | TopK         | 63.57        | 37.73        | 38.04        | 29.77        |
|                  |              | Div          | 37.43        | 23.23        | 26.51        | 18.07        |
|                  |              | TopK-Div     | <b>72.14</b> | <b>44.32</b> | <b>53.99</b> | <b>40.68</b> |
|                  | 8            | Rand         | 24.89        | 18.64        | 21.16        | 9.52         |
|                  |              | TopK         | 69.64        | 44.09        | 56.04        | 41.14        |
| Div              |              | 42.71        | 26.00        | 30.59        | 20.68        |              |
| TopK-Div         | <b>79.29</b> | <b>47.73</b> | <b>64.24</b> | <b>44.32</b> |              |              |
| CodeLlama-34B-hf | 4            | Rand         | 15.75        | 13.02        | 14.42        | 5.98         |
|                  |              | TopK         | 63.57        | 42.05        | 43.51        | 30.23        |
|                  |              | Div          | 39.86        | 24.75        | 32.92        | 19.18        |
|                  |              | TopK-Div     | <b>72.50</b> | <b>48.18</b> | <b>56.72</b> | <b>44.55</b> |
|                  | 8            | Rand         | 25.46        | 20.50        | 24.76        | 11.73        |
|                  |              | TopK         | 73.93        | 48.41        | 56.04        | 44.09        |
| Div              |              | 44.18        | 27.50        | 39.29        | 24.32        |              |
| TopK-Div         | <b>80.71</b> | <b>50.00</b> | <b>64.92</b> | <b>48.86</b> |              |              |

Table 14: Results of GPT-4o-mini and Deepseek-v3 in Math Task.

| Model       | $K$ | Dataset | Method       |              |              |          |
|-------------|-----|---------|--------------|--------------|--------------|----------|
|             |     |         | Rand         | TopK         | Div          | TopK-Div |
| GPT-4o-mini | 4   | GSM8K   | <b>93.03</b> | 91.06        | 92.80        | 92.27    |
|             |     | PRM800K | 68.40        | 66.60        | <b>71.20</b> | 69.20    |
| Deepseek-v3 | 4   | GSM8K   | <b>96.13</b> | 95.75        | 95.91        | 95.45    |
|             |     | PRM800K | 85.00        | <b>87.00</b> | 85.00        | 86.80    |

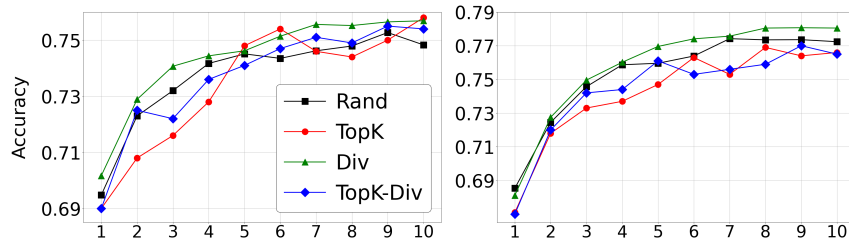


Figure 6: (The accuracy of different methods with different number of shots  $K$  on reading comprehension tasks.) We choose report the results on Llama-3.1-8B, with Sentence-BERT embeddings (all-roberta-large-v1). **Left:** results where demonstration and query come from SCIQ. **Right:** results where demonstration and query come from SQuAD.

We present the experimental results on reading comprehension task (SCIQ, SQuAD), where diversity-aware methods perform well, with different numbers of shots  $K$  ranging from 1 to 10. We test different methods on the Llama-3.1-8B model. Figure 6 summarizes. To our surprise, when  $k = 1$ , TopK performs significantly worse than Rand on both datasets, indicating that the accuracy of these datasets is not solely related to the coverage of example sets. Under most settings of  $k$ , Div shows significant advantages over

1346  
1347  
1348  
1349  
1350

Table 15: Results of K-Means Method in Math Task. We add the K-Means baseline based on Table 1 in our paper. The implementation of K-Means consists of two steps: First, partition the input into  $k$  clusters using the  $k$ -Means method. Second, select  $k$  points as demonstrations by choosing the point closest to the cluster center within each cluster.

| model           | $K$ | Dataset    | Method |       |              |              |              |
|-----------------|-----|------------|--------|-------|--------------|--------------|--------------|
|                 |     |            | Rand   | Topk  | Div          | Topk-Div     | K-means      |
| Llama-3.1-8B    | 4   | GSM8K      | 82.24  | 81.99 | 82.14        | 81.74        | <b>83.89</b> |
|                 |     | GSM-Plus-M | 66.90  | 65.30 | 66.92        | 66.12        | <b>68.10</b> |
|                 | 8   | GSM8K      | 82.81  | 82.26 | <b>82.98</b> | 82.63        | 82.36        |
|                 |     | GSM-Plus-M | 66.72  | 65.99 | 66.56        | 66.48        | 66.52        |
| Gemma-2-9B      | 4   | GSM8K      | 84.29  | 84.25 | 84.44        | 83.54        | <b>85.24</b> |
|                 |     | GSM-Plus-M | 74.40  | 74.50 | 73.34        | 74.47        | 74.52        |
|                 | 8   | GSM8K      | 84.34  | 83.57 | 84.22        | 83.85        | <b>84.97</b> |
|                 |     | GSM-Plus-M | 74.48  | 75.36 | 74.69        | 75.16        | <b>76.29</b> |
| Mistral-7B-v0.3 | 4   | GSM8K      | 48.78  | 49.28 | 49.49        | <b>49.99</b> | 43.90        |
|                 |     | GSM-Plus-M | 37.20  | 38.20 | 37.50        | <b>38.45</b> | 35.50        |
|                 | 8   | GSM8K      | 47.86  | 48.43 | 48.33        | 48.60        | <b>49.13</b> |
|                 |     | GSM-Plus-M | 36.32  | 37.35 | 36.12        | <b>37.81</b> | 36.41        |

Table 16: Embedding on answer using Gemma-2-9B with 4 shots. Comparing to Table 1, the relative ranking between the tested methods doesn’t change.

|                    | Rand  | TopK  | Div   | TopK-Div     |
|--------------------|-------|-------|-------|--------------|
| GSM8K              | 82.21 | 84.53 | 84.14 | <b>84.69</b> |
| PRM800K            | 38.04 | 45.60 | 37.56 | <b>46.40</b> |
| GeoQuery(Standard) | 13.71 | 79.64 | 54.32 | <b>84.29</b> |

TopK. Moreover, the correlation between example sets selected by Div and test samples is relatively low. This sufficiently demonstrates that even when example samples do not have coverage of test samples, they can still be high-quality examples, which is also consistent with the good performance of Rand.

### D.3 More diversity-aware method

In the main text, TopK-Div and Div are both diversity-aware methods that combine the TopK method. We want to understand what happens when using a purely diversity-based method. Therefore, we implemented the K-Means method: dividing the training set into  $k$  clusters by  $k$ -means algorithm and then choose the nearest sample to the Centroid from each cluster (K-Means), K-Means can be viewed as a purely diversity-based met.

The results in Table 15 show that the K-Means method still has advantages compared to the TopK method. In fact, we believe Rand can also be considered a purely diversity-based method. This implies the advantage of diversity methods does not depend on the specific implementation.

### D.4 Abalations on the size of training set

To investigate whether the way diversity works is related to the size of the training set—for example, whether the example selection strategy needs to change when the available training set is limited, We conducted experiments on the SQuAD and SCIQ datasets by randomly sampling 50 examples from each training set to create SCIQ-50 and SQuAD-50, while keeping the original testing set unchanged.

When the available training set size is reduced, TopK still underperforms compared to Div, maintaining an average performance gap of 1% in 4-shot and 8-shot settings.

### D.5 Ablations on “better” embeddings

“better” embedding in a cheating way. All methods we test, except randomly chosen (Rand), depend on an embedding model. It is always possible that the embedding model is not good enough. Indeed, using Sentence-BERT on questions/input (optimized for semantic similarity) might not be optimal for math

Table 17: **Results of different embeddings on Llama-3.1-8B.** We test different methods using different similarity scores computation (“all-roberta-large-v1”, “BM25”, “BertScore”). We test on Llama-3.1-8B model on math (using instruct model) and reading comprehension (using base model). The numbers for embedding “all-roberta-large-v1” are copied from Table 1. The numbers corresponding to Rand for BM25 and BertScore are also copied. We find that: (1) using another embedding might affect the TopK performance, as we can observe an increase of performance for TopK while changing to BM25 or BertScore. (2) Diversity still helps, since if we look at the best performance with the best embedding, in most of the cases the best performance is still achieved by diversity-aware methods.

| Embedding            | K | Method   | Math         |              | Reading      |              |
|----------------------|---|----------|--------------|--------------|--------------|--------------|
|                      |   |          | GSM8K        | PRM800K      | SQuAD        | SCIQ         |
| all-roberta-large-v1 | 4 | Rand     | 82.40        | 43.50        | 75.87        | 74.17        |
|                      |   | TopK     | <b>82.64</b> | 41.40        | 73.70        | 72.80        |
|                      |   | Div      | 82.43        | <b>44.86</b> | <b>76.02</b> | <b>74.44</b> |
|                      |   | TopK-Div | 81.43        | 44.80        | 74.40        | 73.60        |
|                      | 8 | Rand     | 82.77        | 43.32        | 77.35        | 74.79        |
|                      |   | TopK     | 82.11        | 43.00        | 76.90        | 74.40        |
|                      |   | Div      | <b>83.13</b> | <b>44.28</b> | <b>78.05</b> | <b>75.52</b> |
|                      |   | TopK-Div | 81.73        | 40.00        | 75.90        | 74.90        |
|                      |   |          |              |              |              |              |
| BM25                 | 4 | Rand     | 82.40        | 43.50        | 75.87        | 74.17        |
|                      |   | TopK     | 81.88        | 42.00        | 73.80        | <b>74.40</b> |
|                      |   | Div      | <b>82.47</b> | 44.12        | <b>76.65</b> | 72.74        |
|                      |   | TopK-Div | 81.20        | <b>45.00</b> | 75.50        | 74.30        |
|                      | 8 | Rand     | 82.77        | 43.32        | 77.35        | 74.79        |
|                      |   | TopK     | 82.94        | 44.80        | 76.60        | <b>75.20</b> |
|                      |   | Div      | <b>83.44</b> | 43.92        | <b>78.97</b> | 74.12        |
|                      |   | TopK-Div | 83.02        | <b>45.60</b> | 77.50        | 74.50        |
|                      |   |          |              |              |              |              |
| BertScore            | 4 | Rand     | 82.40        | 43.50        | <b>75.87</b> | 74.17        |
|                      |   | TopK     | 81.58        | <b>45.60</b> | 75.00        | <b>74.30</b> |
|                      |   | Div      | <b>82.81</b> | 44.06        | 74.16        | 73.06        |
|                      |   | TopK-Div | 81.05        | 44.40        | 74.90        | 73.20        |
|                      | 8 | Rand     | 82.77        | 43.32        | <b>77.35</b> | 74.79        |
|                      |   | TopK     | 82.34        | 42.60        | 76.40        | <b>75.50</b> |
|                      |   | Div      | <b>83.09</b> | 43.68        | 76.00        | 74.95        |
|                      |   | TopK-Div | 81.58        | <b>44.00</b> | 75.70        | 74.70        |
|                      |   |          |              |              |              |              |

tasks and text-to-SQL generation, and the ideal embedding might be highly dependent on the structure or reasoning steps of the answer. In this section, we test if diversity still helps when given a better embedding, computed in a “cheating” way: For math problems, we append the gold answer after the question and compute the embedding using Sentence-BERT; For text-to-SQL generation, we compute the occurrence of keywords in the answer (Levy et al., 2023). Table 16 summarizes the result using the “cheating” embeddings on Gemma-2-9B, and in general, diversity still helps for these tasks.

**Computing local structure for GeoQuery.** For the code-standard task, we tokenized the sample answers at the word level and obtained 52 distinct tokens, with each dimension representing a token. For a given sample, in its 52-dimensional vector, if the corresponding token appears in its answer, the value at that position is 1, otherwise 0. We use this embedding as the code embedding on answers.

**BM25 and BertScore for math and reading comprehension.** We conduct ablation studies on the model to compute the similarity score, changing from cosine similarity from embeddings computed by “all-roberta-large-v1” to BM25 and BertScore, and test different methods on math and reading comprehension tasks. Table 17 summarizes our results. We find that (1) using another embedding might affect the TopK performance, as we can observe an increase in performance for TopK while changing to BM25 or BertScore. (2) Diversity still helps since if we look at the best performance with the best embedding, in most cases, the best performance is still achieved by diversity-aware methods.

## D.6 Decoding method

In this part we show some preliminary results on changing the decoding strategy for reading comprehension tasks (SQuAD and CommonsenseQA), since for code and math, greedy decoding is known to perform

Table 18: Decode performance using Llama-3.1-8B on reading comprehension tasks. The number of shot is fixed as 4.

| Decode   | Test. | Rand  | TopK  | Div          | TopK-Div     |
|----------|-------|-------|-------|--------------|--------------|
| Greedy   | Squad | 75.87 | 73.70 | <b>76.02</b> | 74.40        |
|          | Sciq  | 74.17 | 72.80 | <b>74.44</b> | 73.60        |
| Sampling | Squad | 70.93 | 72.40 | 70.95        | <b>72.80</b> |
|          | Sciq  | 66.86 | 66.70 | 67.08        | <b>67.70</b> |

well. By changing greedy decoding to sampling decoding (topP = 0.95, Temperature = 0.7), we find that the performance of all tasks drops a lot (Table 18), which justifies our decoding strategy selection.

## E Theoretical justification and simulations

In this section, we give a theoretical justification for combining diversity in demonstration selection for ICL, even if the “embedding” is accurate (Theorems E.2 and E.3). Then we validate the superiority of TopK-Div compared to TopK in more general settings. In Appendix E.3 and Appendix E.4, we also employ the theoretical framework to conduct detailed simulation experiments.

We consider the linear regression model, where there is a task vector  $\theta_{\mathcal{T}} \in \mathbb{R}^d$ . The data for this task has embedded input  $e \in \mathbb{R}^d$  and output  $y = \langle \theta_{\mathcal{T}}, e \rangle$ . We also have a demonstration set  $D = \{(e_i, y_i)\}_{i=1}^n$  with size  $n$ , where  $y_i = \langle \theta_{\mathcal{T}}, e_i \rangle$  and  $e_i$  is drawn from the demonstration distribution  $\mathcal{D}_{\mathcal{E}}$ . Now given a query  $e_q$  drawn from the query distribution  $\mathcal{Q}_{\mathcal{E}}$ , the goal of demonstration selection is to select a subset  $S = \{(e_{j_i}, y_{j_i})\}_{i=1}^K$ , such that given the demonstrations  $S = \{(e_{j_i}, y_{j_i})\}_{i=1}^K$ , the LLM predicts the output close to the gold label  $y_q = \langle \theta_{\mathcal{T}}, e_q \rangle$ , i.e.  $y_q \approx \text{LLM}(S, e_q)$ . We make the following assumption on the mechanism of LLM for learning linear regression in-context, given the demonstration  $S$  and the query  $e_q$ .

**Assumption E.1 (ICL for linear regression).** Suppose that the task  $\mathcal{T}$  is to predict the value of a linear function  $y = \langle \theta_{\mathcal{T}}, e \rangle$  and  $K$  demonstrations  $S = \{(e_{j_i}, y_{j_i})\}_{i=1}^K$  are selected. Denote  $E = [e_{j_1}, \dots, e_{j_K}]^{\top} \in \mathbb{R}^{K \times d}$  as the data matrix. Then given a query  $e_q$ , we assume that the prediction given by the LLM is  $y_{\text{pred}} = \langle e_q, E^{\dagger} E \theta_{\mathcal{T}} \rangle$ . Namely, the LLM learns the min-norm solution for the overparameterized linear regression.

By this assumption, the prediction loss of  $e_q$  is

$$\text{Loss}(e_q) := (y_{\text{pred}} - \langle \theta_{\mathcal{T}}, e_q \rangle)^2 = \langle \theta_{\mathcal{T}} - E^{\dagger} E \theta_{\mathcal{T}}, e_q \rangle^2.$$

ICL for linear regression has been extensively studied, empirically and theoretically (Appendix A). Theorem E.1 is also empirically justified, where (Akyürek et al., 2023) observed that after pretraining an autoregressive transformer model on noiseless linear regression tasks, the transformer will learn the min-norm solution for the linear regression in-context if the size of demonstrations  $K < d$ .

We further assume that the embedding for each data  $e \in \{0, 1\}^d$ . This is inspired by the theoretical framework that each problem from a specific task contains certain skills (or local structures), and an LLM is able to solve that problem perfectly if the LLM knows all the skills (local structures) and is able to compose the skills (local structures) together (Arora and Goyal, 2023; Yu et al., 2024; Zhao et al., 2025). For example, for a specific math problem related to algebra, the skills required to solve this problem are polynomial multiplication and solving equations, while for another math problem related to geometry, the skills required might be changed to coordinate systems and solving equations. It is also worth noting that the skill(local structure)-based embedding design also gains empirical success. For example, (Levy et al., 2023; Didolkar et al., 2024) improves semantic parsing and (An et al., 2023b) improves math ability by selecting demonstrations that require similar skills or local structures to the query.

**Example I: Diversity benefits from coverage.** We characterize the demonstration distribution  $\mathcal{D}_{\mathcal{E}}$  and the query distribution  $\mathcal{Q}_{\mathcal{E}}$  below. Let  $l \geq 200$  be an even number and let  $d = 4l$ , where the choice of 200 is to simplify the analysis. Let  $\mathcal{D}_{\mathcal{E}}$  be: Uniformly draw a subset  $T_1 \subseteq [2l]$  of size  $l/2$  and a subset  $T_2 \subseteq \{2l+1, \dots, 4l\}$  of size  $l/2$ , and output  $e = e_{T_1 \cup T_2}$ , i.e., the  $i$ -th entry of  $e$  is 1 iff  $i \in T_1 \cup T_2$ . Assume the size  $n$  of  $D$  is sufficiently large that  $D$  covers the entire ground set of  $\mathcal{D}_{\mathcal{E}}$ . Let  $\mathcal{Q}_{\mathcal{E}}$  be: Uniformly draw a subset  $T \subset [2l]$  of size  $l$ . We have the following theorem, whose proof can be found in Appendix E.1.

**Theorem E.2 (Justification example I).** Suppose each entry of  $\theta_{\mathcal{T}}$  is i.i.d. drawn from the uniform distribution on  $[0, 1]$ . Let  $K = 2$  and  $\mathcal{D}_{\mathcal{E}}, \mathcal{Q}_{\mathcal{E}}$  be as defined above. For a query  $e_q$  drawn from  $\mathcal{Q}_{\mathcal{E}}$ , let  $L, L'$  denote the expected prediction loss of  $e_q$  using  $\text{TopK}$  and  $\text{TopK-Div}$ , respectively, where the randomness comes from  $\theta_{\mathcal{T}}, e_q$ , and the selection of demonstration examples. Then  $L > L'$  for any hyperparameter  $\alpha \in (0, 1)$  for  $\text{TopK-Div}$ .

Intuitively, the selected two demonstration examples of  $\text{TopK-Div}$  must cover all non-zero entries of  $e_q$ , while this property is unlikely to hold for  $\text{TopK}$ . This demonstrates that adding diversity may increase the coverage of demonstration examples to queries and lead to a lower prediction loss, aligning with the findings in (Levy et al., 2023; Gupta et al., 2023; Ye et al., 2023).

**Example II: Diversity is beyond coverage.** We again characterize  $\mathcal{D}_{\mathcal{E}}$  and  $\mathcal{Q}_{\mathcal{E}}$  below. Let  $l \geq 3$  be an integer and let  $d = 4l$ . Let  $\mathcal{D}_{\mathcal{E}}$  be: Uniformly draw a subset  $T_1 \subseteq [2l]$  of size  $l - 1$  and a subset  $T_2 \subseteq \{2l + 1, \dots, 4l\}$  of size 1, and output  $e = e_{T_1 \cup T_2}$ . Assume the size  $n$  of  $D$  is sufficiently large that  $D$  covers the entire ground set of  $\mathcal{D}_{\mathcal{E}}$ . Let  $\mathcal{Q}_{\mathcal{E}}$  be: Uniformly draw a subset  $T \subset [2l]$  of size  $l$ . We have the following theorem for this example, whose proof can be found in Appendix E.2.

**Theorem E.3 (Justification example II).** Suppose each entry of  $\theta_{\mathcal{T}}$  is i.i.d. drawn from the uniform distribution on  $[0, 1]$ . Let  $K = 2$  and  $\mathcal{D}_{\mathcal{E}}, \mathcal{Q}_{\mathcal{E}}$  be as defined above. For a query  $e_q$  drawn from  $\mathcal{Q}_{\mathcal{E}}$ , let  $L, L'$  denote the expected prediction loss of  $e_q$  using  $\text{TopK}$  and  $\text{TopK-Div}$ , respectively, where the randomness comes from  $\theta_{\mathcal{T}}, e_q$ , and the selection of demonstration examples. Then  $L > L'$  if hyperparameter  $\alpha \geq 1 - 1/l$  for  $\text{TopK-Div}$ .

The demonstration examples of  $\text{TopK}$  and  $\text{TopK-Div}$  must cover all non-zero entries of  $e_q$ . The smaller loss of  $\text{TopK-Div}$  is caused by selecting two demonstration examples with different non-zero entries among  $\{2l + 1, \dots, 4l\}$ , indicating that adding diversity could benefit ICL “beyond coverage”.

In Appendix E.3, we conduct simulations to validate that the advantage of  $\text{TopK-Div}$  over  $\text{TopK}$ , driven by coverage and beyond, extends to more general settings, including the ID setting ( $\mathcal{D}_{\mathcal{E}} = \mathcal{Q}_{\mathcal{E}}$ ) and scenarios with different training scales for  $D$ .

## E.1 Proof of Theorem E.2: justification example I

Fix a query  $e_q$  drawn from  $\mathcal{Q}_{\mathcal{E}}$ . By symmetry, we can assume the non-zero entry set of  $e_q$  is  $[2l]$ . For simplicity, we let  $\theta = \theta_{\mathcal{T}}$ .

**Demonstration example set for TopK-Div** We first analyze the demonstration example set for  $\text{TopK-Div}$ , denoted by  $S = \{s^{(1)}, s^{(2)}\} \subseteq D$ . Let  $T^{(t)}$  denote the non-zero entry set of  $s^{(t)}$ . By the construction of  $\mathcal{D}$ , we first note that  $|T^{(1)} \cap [l]| = \frac{l}{2}$ . By the rule of  $\text{TopK-Div}$ , we also note that  $|T^{(2)} \cap [l]| = \frac{l}{2}$  and  $T^{(1)} \cap T^{(2)} = \emptyset$ . Such  $s^{(2)}$  must exist since all elements in the ground set of  $\mathcal{D}_{\mathcal{E}}$  are contained in  $D$ , and is selected since it minimizes

$$\alpha \cdot \text{Similarity}(e, e_q) + (1 - \alpha) \text{Diversity}(e, S)$$

over all  $e \in D - \{s^{(1)}\}$ .

**Demonstration example set for TopK** Next, we compute the expected prediction loss  $L$  for  $\text{TopK}$ . Again, let its demonstration example set be  $S = \{s^{(1)}, s^{(2)}\} \subseteq D$ . Let  $T^{(t)}$  denote the non-zero entry set of  $s^{(t)}$ . By the construction of  $\mathcal{D}$ , we note that  $|T^{(1)} \cap [l]| = |T^{(2)} \cap [l]| = \frac{l}{2}$ . However, different from the case of  $\text{TopK-Div}$ ,  $|T^{(1)} \cap T^{(2)}|$  can vary from 0 to  $l - 1$ . To handle this, we define  $a = |T^{(1)} \cap T^{(2)} \cap [l]|$  and  $b = |T^{(1)} \cap T^{(2)} \cap ([d] \setminus [l])|$ , and define  $L_{a,b}$  to be the expected prediction loss conditioned on pair  $(a, b)$ . Note that  $0 \leq a, b \leq l/2$  and  $a + b \leq l - 1$ .

**Comparing  $L$  and  $L'$**  We remark that  $L$  is a linear combination  $\sum_{a,b} p_{a,b} L_{a,b}$  with  $\sum_{a,b} p_{a,b} = 1$ , where  $p_{a,b}$  is the conditional probability with respect to intersection numbers  $(a, b)$ . Also,  $L' = L_{0,0}$ . By symmetry, we have the following observation:

$$\Pr[a \leq l/4 \leq b] \geq 0.25,$$

where  $l/4$  is the expectation of  $a$  and  $b$ . Thus, we have

$$L \geq \sum_{a \leq l/4 \leq b} p_{a,b} L_{a,b} \geq \sum_{a,b \in l/4 \pm \sqrt{l}} p_{a,b} \cdot \min_{a \leq l/4 \leq b} L_{a,b} \geq 0.25 \min_{a \leq l/4 \leq b} L_{a,b}.$$

Thus, to prove  $L > L'$ , it suffices to prove the following lemma.

**Lemma E.4 (Comparing  $L_{a,b}$  and  $L_{0,0}$ ).** *For any  $a \leq l/4 \leq b$ , we have  $L_{a,b} > 4L_{0,0}$ .*

*Proof.* By symmetry, we assume  $T^{(1)} = [\frac{l}{2}] \cup ([\frac{5}{2}l] - [2l])$ ,  $T^{(2)} = ([l] - [a] - [\frac{l}{2} - a]) \cup ([3l - b] - [\frac{5}{2}l - b])$ ,  $|T^{(1)} \cap T^{(2)} \cap [L]| = |T^{(1)} \cap T^{(2)} \cap [2L]| = a$ ,  $|T^{(1)} \cap T^{(2)} \cap ([4L] - [2L])| = b$ . The expected prediction loss for this setting equals  $L_{a,b}$  since  $\theta_i$ s are i.i.d. random variables. Let  $\theta$  denote the min-norm solution defined as in Assumption E.1. Then we have

$$\langle \hat{\theta} - \theta, e_{T^{(1)}} \rangle = \sum_{i=1}^{\frac{l}{2}} \hat{\theta}_i + \sum_{i=2l+1}^{\frac{5}{2}l} \hat{\theta}_i - \sum_{i=1}^{\frac{l}{2}} \theta_i - \sum_{i=2l+1}^{\frac{5}{2}l} \theta_i = 0, \quad (5)$$

and

$$\langle \hat{\theta} - \theta, e_{T^{(2)}} \rangle = \sum_{i=\frac{l}{2}-a+1}^{l-a} \hat{\theta}_i + \sum_{i=\frac{5}{2}l-b+1}^{3l-b} \hat{\theta}_i - \sum_{i=\frac{l}{2}-a+1}^{l-a} \theta_i - \sum_{i=\frac{5}{2}l-b+1}^{3l-b} \theta_i = 0. \quad (6)$$

To get the min-norm solution, we need to minimize the following Lagrangian multiplier

$$\mathcal{L}(\hat{\theta}, \lambda_1, \lambda_2) = \sum_{i=1}^{l-a} \hat{\theta}_i^2 - 2\lambda_1 \langle \hat{\theta} - \theta, e_{T^{(1)}} \rangle - 2\lambda_2 \langle \hat{\theta} - \theta, e_{T^{(2)}} \rangle.$$

To ensure the partial derivatives with respect to  $\hat{\theta}$  equal to 0, we obtain that

$$\begin{aligned} \hat{\theta}_1 &= \dots = \hat{\theta}_{\frac{l}{2}-a} = \hat{\theta}_{2l+1} = \dots = \hat{\theta}_{\frac{5}{2}l-b} = \lambda_1, \\ \hat{\theta}_{\frac{l}{2}+1} &= \dots = \hat{\theta}_{l-a} = \hat{\theta}_{\frac{5}{2}l+1} = \dots = \hat{\theta}_{3l-b} = \lambda_2, \\ \hat{\theta}_{\frac{l}{2}-a+1} &= \dots = \hat{\theta}_{\frac{l}{2}} = \hat{\theta}_{\frac{5}{2}l-b+1} = \dots = \hat{\theta}_{\frac{5}{2}l} = \lambda_1 + \lambda_2, \\ \hat{\theta}_{l-a+1} &= \dots = \hat{\theta}_{2l} = \hat{\theta}_{3l-b+1} = \dots = \hat{\theta}_{4l} = 0. \end{aligned} \quad (7)$$

Adding Equations (5)-(7), we have

$$(l + a + b)(\lambda_1 + \lambda_2) = \sum_{i=1}^{\frac{l}{2}} \theta_i + \sum_{i=2l+1}^{\frac{5}{2}l} \theta_i + \sum_{i=\frac{l}{2}-a+1}^{l-a} \theta_i + \sum_{i=\frac{5}{2}l-b}^{3l-b} \theta_i. \quad (8)$$

Thus, we conclude that

1498

$$\begin{aligned}
& \left[ \sum_{i=1}^l \widehat{\theta}_i - \sum_{i=1}^l \theta_i \right]^2 \\
&= \left[ \left( \frac{l}{2} - a \right) (\lambda_1 + \lambda_2) + a\lambda_1 + a\lambda_2 - \sum_{i=1}^l \theta_i \right]^2 \\
&= \left[ \frac{l}{2} (\lambda_1 + \lambda_2) - \sum_{i=1}^l \theta_i \right]^2 \\
&= \left[ \frac{\frac{l}{2}}{l+a+b} \left( \sum_{i=1}^{\frac{l}{2}} \theta_i + \sum_{i=2l+1}^{\frac{5l}{2}} \theta_i + \sum_{i=\frac{l}{2}-a+1}^{l-a} \theta_i + \sum_{i=\frac{5l}{2}-b}^{3l-b} \theta_i \right) - \sum_{i=1}^l \theta_i \right]^2 \\
&= \left[ -\frac{\frac{l}{2} + a + b}{l+a+b} \sum_{i=1}^{\frac{l}{2}-a} \theta_i - \frac{\frac{l}{2} + a + b}{l+a+b} \sum_{i=\frac{l}{2}+1}^{l-a} \theta_i - \frac{a+b}{l+a+b} \sum_{i=\frac{l}{2}-a+1}^{\frac{l}{2}} \theta_i \right. \\
&\quad \left. + \frac{\frac{l}{2}}{l+a+b} \sum_{i=2l+1}^{\frac{5l}{2}} \theta_i + \frac{\frac{l}{2}}{l+a+b} \sum_{i=\frac{5l}{2}-b+1}^{3l-b} \theta_i \right]^2,
\end{aligned}$$

1499

where the first equation follows from Equation (7) and the third equation follows from Equation (8). Since each  $\theta_i$  is i.i.d. drawn from the uniform distribution over  $[0, 1]$ , we have

1500

1501

$$\begin{aligned}
L_{a,b} &= \mathbb{E} \left[ \langle \widehat{\theta} - \theta, e_q \rangle^2 \right] \\
&= \mathbb{E} \left[ \left[ \sum_{i=1}^l \widehat{\theta}_i - \sum_{i=1}^l \theta_i \right]^2 \right] \\
&= \frac{\left( \frac{l}{2} + a + b \right)^2 \left( \frac{l}{2} - a \right) + \frac{a(a+b)^2}{2} + \frac{l^3}{8} + \frac{3(bl-a^2-ab)^2}{2}}{6(l+a+b)^2}.
\end{aligned}$$

1502

1503

1504

Thus,  $L_{0,0} = \frac{l}{24}$ . When  $a \leq l/4 \leq b$ , we have

1505

$$\begin{aligned}
L_{a,b} &> \frac{3(bl - a^2 - ab)^2/2}{6(l+a+b)^2} \\
&\geq \frac{(l^2/4 - 2(l/4)^2)^2}{4(2l)^2} && (a \leq l/4 \leq b) \\
&= \frac{(l^2/8)^2}{16l^2} \\
&= \frac{l^2}{1024} \\
&\geq 4L_{0,0}. && (l \geq 200)
\end{aligned}$$

1506

1507

1508

1509

1510

This completes the proof.  $\square$

1511

## E.2 Proof of Theorem E.3: justification example II

1512

By symmetric, we fix  $e_q = e_{[l]}$ . Like the proof of Theorem E.2, we first study the demonstration example sets, denoted by  $S = \{s^{(1)}, s^{(2)}\} \subseteq D$ , derived from TopK and TopK-Div. We observe that for both

1513

1514

algorithms,  $|T^{(1)} \cap [l]| = |T^{(2)} \cap [l]| = l - 1$ . Note that this property for  $\text{TopK-Div}$  follows from the choice of  $\alpha \geq 1 - \frac{1}{l}$ , which ensures that  $|T^{(2)} \cap [l]| \leq l - 2$  can not achieve the minimum for

$$\alpha \cdot \text{Similarity}(e, e_q) + (1 - \alpha) \text{Diversity}(e, S)$$

Thus, by symmetry, we can fix  $T^{(1)} = [l - 1] \cup \{2l + 1\}$  and there are only three choices for  $T^{(2)}$ :

- Case 1:  $T^{(2)} = [l] \cup \{2l + 2\} - \{1\}$ ;
- Case 2:  $T^{(2)} = [l] \cup \{2l + 1\} - \{1\}$ .
- Case 3:  $T^{(2)} = [l - 1] \cup \{2l + 2\}$ .

We define the expected prediction loss of these three cases to be  $L_1, L_2, L_3$ , respectively. By the definition of  $\text{TopK-Div}$ , we know that  $L' = L_1$ . Moreover, the expected prediction loss  $L$  of  $\text{TopK}$  must be a linear combination of  $L_1, L_2, L_3$ . Thus, it suffices to prove that  $L_2 > L_1$  and  $L_3 > L_1$ . Below, we compute  $L_1, L_2, L_3$  separately.

**Computing  $L_1$ .** The computation idea is similar to that of Lemma E.4. Suppose  $\hat{\theta}$  is the min-norm solution and we have

$$\sum_{i=1}^{l-1} \hat{\theta}_i + \hat{\theta}_{2l+1} - \sum_{i=1}^{l-1} \theta_i - \theta_{2l+1} = 0 \text{ and } \sum_{i=2}^l \hat{\theta}_i + \hat{\theta}_{2l+2} - \sum_{i=2}^l \theta_i - \theta_{2l+2} = 0.$$

Again, consider the Lagrangian multiplier  $\mathcal{L}(\hat{\theta}, \lambda_1, \lambda_2) = \sum_{i=1}^l (\hat{\theta}_i)^2 - 2\lambda_1 \langle \hat{\theta} - \theta, e_{T^{(1)}} \rangle - 2\lambda_2 \langle \hat{\theta} - \theta, e_{T^{(2)}} \rangle$ .

To ensure the partial derivative w.r.t.  $\hat{\theta}$  equal to 0, we have

$$\hat{\theta}_2 = \hat{\theta}_3 = \dots = \hat{\theta}_{l-1} = \lambda_1 + \lambda_2, \text{ and } \hat{\theta}_1 = \hat{\theta}_{2l+1} = \lambda_1, \hat{\theta}_l = \hat{\theta}_{2l+2} = \lambda_2.$$

Combining the above equations, we have

$$(2l - 2)(\lambda_1 + \lambda_2) = 2 \sum_{i=2}^{l-1} \theta_i + \theta_1 + \theta_l + \theta_{2l+1} + \theta_{2l+2}.$$

Thus,

$$\left( \sum_{i=1}^l \theta_i - \sum_{i=1}^l \hat{\theta}_i \right)^2 = [(l - 1)(\lambda_1 + \lambda_2) - \sum_{i=1}^l \theta_i]^2 = \left( \frac{\theta_{2l+1} + \theta_{2l+2} - \theta_1 - \theta_l}{2} \right)^2.$$

Consequently, we have

$$L_1 = \mathbb{E}[\langle \hat{\theta} - \theta, e_q \rangle^2] = \mathbb{E}[\left( \frac{\theta_{2l+1} + \theta_{2l+2} - \theta_1 - \theta_l}{2} \right)^2] = \frac{1}{12}.$$

**Computing  $L_2$ .** Similarly, we have

$$\sum_{i=1}^{l-1} \hat{\theta}_i + \hat{\theta}_{2l+1} - \sum_{i=1}^{l-1} \theta_i - \theta_{2l+1} = 0, \text{ and } \sum_{i=2}^l \hat{\theta}_i + \hat{\theta}_{2l+1} - \sum_{i=2}^l \theta_i - \theta_{2l+1} = 0.$$

Thus, using the Lagrangian multiplier, we obtain that

$$\hat{\theta}_2 = \hat{\theta}_3 = \dots = \hat{\theta}_{l-1} = \hat{\theta}_{2l+1} = \lambda_1 + \lambda_2, \text{ and } \hat{\theta}_1 = \lambda_1, \hat{\theta}_l = \lambda_2.$$

Combining the above equations, we have

$$(2l - 1)(\lambda_1 + \lambda_2) = 2 \sum_{i=2}^{l-1} \theta_i + 2\theta_{2l+1} + \theta_1 + \theta_l.$$

Thus,

$$\begin{aligned} L_2 &= \mathbb{E}\left[\left(\sum_{i=1}^l \theta_i - \sum_{i=1}^l \widehat{\theta}_i\right)^2\right] = \mathbb{E}\left[\sum_{i=1}^l \theta_i - [(l-1)(\lambda_1 + \lambda_2)]\right]^2 \\ &= \mathbb{E}\left[\frac{1}{2l-1} \sum_{i=2}^{l-1} \theta_i + \frac{l}{2l-1}(\theta_1 + \theta_l) - \frac{2l-2}{2l-1}\theta_{2l+1}\right]^2 = \frac{9l^2 - 7l + 2}{12(12l-1)^2} > L_1. \end{aligned}$$

**Computing  $L_3$ .** Similarly, we have

$$\sum_{i=1}^{l-1} \widehat{\theta}_i + \widehat{\theta}_{2l+1} - \sum_{i=1}^{l-1} \theta_i - \theta_{2l+1} = 0, \text{ and } \sum_{i=1}^{l-1} \widehat{\theta}_i + \widehat{\theta}_{2l+2} - \sum_{i=1}^{l-1} \theta_i - \theta_{2l+2} = 0.$$

Using the Lagrangian multiplier, we obtain that

$$\widehat{\theta}_1 = \widehat{\theta}_2 = \widehat{\theta}_3 = \dots = \widehat{\theta}_{l-1} = \lambda_1 + \lambda_2, \text{ and } \widehat{\theta}_{2l+1} = \lambda_1, \widehat{\theta}_{2l+2} = \lambda_2.$$

Combining the above equations, we have

$$(2l-1)(\lambda_1 + \lambda_2) = 2 \sum_{i=1}^{l-1} \theta_i + \theta_{2l+1} + \theta_{2l+2}.$$

Thus,

$$\begin{aligned} L_3 &= \mathbb{E}\left[\left(\sum_{i=1}^l \theta_i - \sum_{i=1}^l \widehat{\theta}_i\right)^2\right] = \mathbb{E}\left[\sum_{i=1}^l \theta_i - [(l-1)(\lambda_1 + \lambda_2)]\right]^2 \\ &= \mathbb{E}\left[\frac{1}{2l-1} \sum_{i=1}^{l-1} \theta_i + \theta_l - \frac{l-1}{2l-1}(\theta_{2l+1} + \theta_{2l+2})\right]^2 = \frac{9l^2 - 7l + 2}{12(2l-1)^2} > L_1. \end{aligned}$$

Overall, we complete the proof of Theorem E.3.

### E.3 Experiment settings

We consider the ID setting with  $\mathcal{D}_{\mathcal{E}} = \mathcal{Q}_{\mathcal{E}}$ .

**Metric for coverage.** Given a sample  $(e, y_E)$ , let  $T^{(e)}$  denote the non-zero entry set of  $e$ . Given a demonstration example set  $S \subseteq D$  and a query  $e_q$ , we define the coverage ratio of  $S$  with respect to  $e_q$  to be:

$$r_S(e_q) := \frac{|\left(\bigcup_{e \in S} T^{(e)}\right) \cap T^{(e_q)}|}{|T^{(e_q)}|},$$

i.e., the ratio of non-zero entries of  $e_q$  covered by samples in  $S$ . By definition,  $r_S(e_q) \in [0, 1]$  and a larger  $r_S(e_q)$  represents higher coverage. Specifically, when  $r_S(e_q) = 1$ , we say  $e_q$  is fully covered by  $S$ . Moreover, given a method  $\mathcal{A}$  that generates a demonstration example set  $A(e_q) \subseteq D$  for each query  $e_q$ , we define

$$r(\mathcal{A}) := \mathbb{E}_{e_q \sim \mathcal{Q}_{\mathcal{E}}}[r_{\mathcal{A}(e_q)}(e_q)] \quad (9)$$

to be the expected value of its coverage ratio  $r_{\mathcal{A}(e_q)}(e_q)$ . If  $r(\mathcal{A}) = 1$ , we say every query is fully covered by  $\mathcal{A}$ .

We want to study the loss difference between TopK-Div and TopK under two scenarios: 1) when query  $e_q$  is fully covered by both algorithms TopK-Div and TopK, i.e.,  $r(\text{TopK-Div}) = r(\text{TopK}) = 1$ ; and 2) when the coverage ratio of TopK-Div is smaller than that of TopK, i.e.,  $r(\text{TopK-Div}) < r(\text{TopK})$ .

Table 19: **(Simulation of the min-norm solution)** “Coverage” represents the coverage ratio of methods, defined as in Equation (9). For each random seed, we selected one hundred test samples. We report the average results across 3 different random seeds for each metric.

| Method   | Shot    | Metric   | Train scale = 1 |         |         | Train scale = 5 |         |         | Train scale = 10 |         |         |
|----------|---------|----------|-----------------|---------|---------|-----------------|---------|---------|------------------|---------|---------|
|          |         |          | $l = 3$         | $l = 4$ | $l = 8$ | $l = 3$         | $l = 4$ | $l = 8$ | $l = 3$          | $l = 4$ | $l = 8$ |
| TopK     | $K = 4$ | Loss     | 0.21            | 0.31    | 12.70   | 0.15            | 0.30    | 9.55    | 0.19             | 0.30    | 7.51    |
|          |         | Coverage | 1.00            | 1.00    | 0.55    | 1.00            | 1.00    | 0.61    | 1.00             | 1.00    | 0.66    |
|          | $K = 8$ | Loss     | 0.47            | 0.57    | 3.09    | 0.43            | 0.84    | 1.33    | 0.45             | 0.83    | 1.19    |
|          |         | Coverage | 1.00            | 1.00    | 0.75    | 1.00            | 1.00    | 0.80    | 1.00             | 1.00    | 0.81    |
| TopK-Div | $K = 4$ | Loss     | 0.19            | 0.32    | 10.25   | 0.18            | 0.31    | 5.47    | 0.21             | 0.29    | 3.97    |
|          |         | Coverage | 1.00            | 1.00    | 0.63    | 1.00            | 1.00    | 0.75    | 1.00             | 1.00    | 0.80    |
|          | $K = 8$ | Loss     | 0.31            | 0.38    | 2.58    | 0.23            | 0.38    | 1.32    | 0.20             | 0.38    | 1.75    |
|          |         | Coverage | 1.00            | 1.00    | 0.87    | 1.00            | 1.00    | 0.94    | 1.00             | 1.00    | 0.94    |

**Parameters.** Let  $d = 200$ . Let  $l$  vary from 3, 4, 8. Let  $K = 4$  or 8. Let  $\mathcal{D}_{\mathcal{E}} = \mathcal{Q}_{\mathcal{E}}$  be the distribution that first samples a subset  $T \subset [d]$  of size  $l$  and then generate  $e_T$ . We set the size of training set  $D$  to be  $|D| = d \times \text{train\_scale}$ , where  $\text{train\_scale} \in \{1, 5, 10\}$ .

For each pair  $(l, K)$ , we generate a testing set  $D_{\text{test}}$  of size 100. We ensure that  $D_{\text{test}} \cap D = \emptyset$ . We report the expected prediction loss and coverage ratio of TopK and TopK-Div for each pair  $(l, K)$ .

#### E.4 Result and discussions

The results, reported in Table 19, reveal key insights into the performance differences between TopK and TopK-Div. We observe that when  $l = 8$ , the coverage ratio of TopK is lower than that of TopK-Div, while its loss is significantly higher. For example, when  $l = 8$ ,  $K = 4$ , and  $\text{train\_scale} = 5$ , the coverage ratio is  $r(\text{TopK}) = 0.61$ , compared to  $r(\text{TopK-Div}) = 0.75$ , while the loss for TopK is 9.55, notably larger than the 5.47 observed for TopK-Div. This demonstrates that incorporating diversity can reduce prediction loss by improving coverage, aligning with Theorem E.2.

When  $l = 3$  or 4, the coverage ratios of TopK and TopK-Div are both 1. We find that the loss of TopK is comparable to or even lower than that of TopK-Div when  $K = 4$ , but significantly higher when  $K = 8$ , across various training scales. For instance, when  $l = 3$ ,  $K = 8$ , and  $\text{train\_scale} = 5$ , the loss for TopK is 0.43, whereas for TopK-Div it is 0.23. This supports our findings in Theorem E.3, demonstrating that diversity can enhance in-context learning beyond just coverage. The inverse trend in loss between  $K = 4$  and  $K = 8$  suggests that increasing coverage is beneficial when the query is not fully covered but becomes redundant when the demonstration example set already provides sufficient coverage.

#### The Use of Large Language Models (LLMs)

We employed large language models (LLMs) solely for word-level grammar checking and minor stylistic refinement of the manuscript. Beyond this limited function, LLMs did not contribute to any other aspects of our research or writing, including conceptualization, experimental design, data analysis, or interpretation of results.