

# Large Language Model-Aware In-Context Learning for Code Generation

Jia Li

Key Lab of HCST (PKU), MOE, SCS,  
Peking University  
Beijing, China  
lijiaa@pku.edu.cn

Ge Li

Key Lab of HCST (PKU), MOE, SCS,  
Peking University  
Beijing, China  
lige@pku.edu.cn

Chongyang Tao

Peking University  
Beijing, China  
chongyangtao@gmail.com

Jia Li ♂, Huangzhao Zhang

Key Lab of HCST (PKU), MOE, SCS,  
Peking University  
Beijing, China  
lijia@stu.pku.edu.cn  
zhang\_hz@pku.edu.cn

Fang Liu

Beihang University  
Beijing, China  
fangliu@buaa.edu.cn

Zhi Jin

Key Lab of HCST (PKU), MOE, SCS,  
Peking University  
Beijing, China  
zhijin@pku.edu.cn

## ABSTRACT

Large language models (LLMs) have shown impressive in-context learning (ICL) ability in code generation. LLMs take a prompt consisting of requirement-code examples and a new requirement as input, and output new programs without any parameter updates. Existing studies have found that the performance of ICL is highly dominated by the quality of selected examples and thus arises research on example selection: given a new requirement, several examples are retrieved from a candidate pool for ICL. However, existing approaches are mostly based on heuristics. They randomly select examples or only consider the textual similarity of requirements to retrieve, leading to sub-optimal performance. In this paper, we propose a novel learning-based selection approach named LAIL (LLM-Aware In-context Learning) for code generation. Given a candidate example, we exploit LLMs themselves to estimate it by considering the generation probabilities of ground-truth programs given a requirement and the example. We then label candidate examples as positive or negative through the probability feedback. Based on the labeled data, we import a contrastive learning objective to train an effective retriever that acquires the preference of LLMs in code generation. We argue that considering the feedback from LLMs themselves is logical and our approach can select suitable examples for ICL since it is consistent with the fact that LLMs have different prior knowledge and preferences. To evaluate our approach, we apply LAIL to three LLMs and evaluate it on three representative datasets (*e.g.*, MBJP, MBPP, and MBCPP). Extensive experiments demonstrate that LATA outperforms the state-of-the-art baselines by 11.58%, 6.89%, and 5.07% on CodeGen, and 4.38%, 2.85%, and 2.74% on GPT-3.5 in terms of Pass@1, respectively. Human evaluation further verifies the superiority of our approach in

three aspects. We also find our LAIL has surprising transferability across LLMs and datasets.

## KEYWORDS

Code generation, in-context-learning, large language model

### ACM Reference Format:

Jia Li, Ge Li, Chongyang Tao, Jia Li ♂, Huangzhao Zhang, Fang Liu, and Zhi Jin. 2018. Large Language Model-Aware In-Context Learning for Code Generation. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Code generation aims to automatically generate the source code given a natural language requirement [7, 18, 30]. It plays an important role in improving the productivity of software development. With the emergence of large language models (LLMs), in-context learning (ICL) [11] becomes a prevailing paradigm and achieves impressive performance in code generation [16, 31, 32]. Given limited examples as the prompt, ICL imitates the human ability to leverage prior knowledge to generate programs without parameter updates. However, ICL comes along with the robustness problem [33], which is sensitive to the selected examples in prompts, resulting in the performance usually varying from almost random to near state-of-the-art performance.

Despite the importance of selecting examples, only a few studies attempt to investigate example selection in code generation [10, 14, 31]. One line of work is to randomly select examples from a candidate pool [14, 26]. The selected examples are usually improper to the test requirement and their semantic distributions vary a lot, resulting in unstable performance. Another line is based on learning-free heuristics [31], which leverages off-the-shelf retrievers such as BM25 [42] to select examples. They calculate the textual similarity between a test requirement and the requirement of examples, then select candidate items with high similarities. Despite the improved performance, these approaches only consider textual matching among requirements and thus are easily biased by exterior lexicon features. Figure 1 reports the top-3 examples by random selection and the BM25 approach for the test requirement “write a function to check whether the entered number is greater than the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference'17, July 2017, Washington, DC, USA*  
© 2018 Association for Computing Machinery.  
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00  
<https://doi.org/XXXXXXXX.XXXXXXX>

**Test Requirement: Write a function to check whether the entered number is greater than the elements of the given array.**

" Note that: To do this, LLMs need to **traverse** the **array** "

- |   |   |   |
|---|---|---|
| <ul style="list-style-type: none"> <li>➤ Write a function to convert camel case string to snake case string by using regex.</li> <li>➤ Write a function to calculate volume of a tetrahedron.</li> <li>➤ Write a Java function to find the sum of fourth power of first n even natural number.</li> </ul> <p>(a) Random TOP-3</p> | <ul style="list-style-type: none"> <li>➤ <u>Write a function to check whether the given number is <b>armstrong</b> or not.</u></li> <li>➤ <u>Write a function to find the <b>lcm</b> of the given array elements.</u></li> <li>➤ <u>Write a Java function to check whether two given <b>lines are parallel</b> or not.</u></li> </ul> <p>(b) BM25 TOP-3</p> | <ul style="list-style-type: none"> <li>➤ Write a function to check if <b>any list element is present</b> in the given list.</li> <li>➤ Write a function to find the <b>cumulative sum of all the values</b> that are present in the tuple list.</li> <li>➤ Write a function to find the difference between largest and smallest value <b>in a given array</b>.</li> </ul> <p>(c) LAIL TOP-3</p> |
|---|---|---|

**Figure 1: Exhibition of the selected top-3 examples by random, BM25, and our LAIL approaches.**

elements of the given array” in MBPP dataset [9]. To accomplish this requirement, LLMs need to perform traversal operations in an array. Randomly selected examples are irrelevant to the test requirement. Although retrieved examples by BM25 have a large amount of textual overlap with the test item labeled by underline, the overlapping words are trivial such as “write a function”. Meanwhile, the operations of non-overlapping words in red color, such as “lines are parallel or not”, can provide less information to LLMs. Examples in ICL are supposed to provide information that LLMs require, and then guide them to generate programs, but the existing approaches leave a huge gap. They do not consider the prior knowledge in LLMs and ignore their preferences. Although one may remedy it by inputting more examples in the prompt, it is highly impractical due to the limited input length of LLMs. Therefore, selecting applicable examples for ICL becomes more urgent.

In this paper, we propose a novel learning-based approach dubbed LAIL to select in-context examples. Instead of selecting examples via textual similarity, given a candidate example, our approach leverages LLMs themselves (LLM-aware) to measure it by incorporating the prediction probability of ground truth given a requirement and the candidate item. To quantify the probability feedback, we design a new metric and label candidate examples based on their metric scores. We treat examples with higher metric scores as positive examples, meanwhile, label examples equipped with lower scores as negative examples since good examples should be beneficial for LLMs to generate correct programs. Based on labeled data, we train a neural retriever through a contrastive learning objective to acquire the preference of LLMs in code generation. From Figure 1, we can find that although not textually similar to the test requirement, the selected examples of LAIL contain the “traversal” operation and are related to the array as shown in purple color tokens.

We evaluate LAIL on three representative LLMs including ChatGPT [1], GPT-3.5 [2], and CodeGen (16B) [47]. We conduct extensive experiments on three datasets (*i.e.*, MBJP (Java) [8], MBPP (Python) [9], and MBCPP (C++) [8]). We use a widely used evaluation metric Pass@k (k = 1, 3, 5) to measure the performance of different approaches. We obtain some findings from experimental results. ❶ In terms of Pass@1, LAIL significantly outperforms the state-of-the-art (SOTA) baselines by 11.58%, 6.89%, and 5.07% on CodeGen, and 4.38%, 2.85%, and 2.74% on GPT-3.5, respectively. ❷ We conduct a human evaluation to measure generated programs

in three aspects (*e.g.* code correctness, quality, and maintainability) and further prove the superiority of LAIL. ❸ Our approach has surprising transferability across LLMs and datasets and brings satisfying improvements to them. ❹ We investigate two plausible choices to estimate candidate examples by LLMs themselves and demonstrate the effectiveness of our design.

We summarize our contributions in this paper as follows:

- Our paper investigates example selection for ICL and argues that a good approach should consider what knowledge LLMs themselves require.
- We propose a novel learning-based approach dubbed LAIL. Given an example, LAIL exploits LLMs themselves to estimate it by calculating the generation probabilities of ground truth given a requirement and the example. We then label candidate examples via their probabilities. Based on labeled data, we import contrastive learning to optimize a retriever.
- We evaluate our approach on three LLMs and conduct extensive experiments on three datasets. Qualitative and quantitative experiments reveal that LAIL significantly outperforms the state-of-the-art baselines.

## 2 BACKGROUND

### 2.1 Large Language Models

In this section, we focus on large language models (LLMs) for source code. LLMs are neural networks that aim to learn the statistical patterns in programming languages [6]. LLMs are pre-trained on a large-scale unlabeled code corpus with the next tokens prediction objective. Specifically, given a program with the token sequence  $C = \{c_1, c_2, \dots, c_n\}$ , LLMs are trained to predict the next token based on some previous tokens:

$$\mathcal{L}_{NTP}(C) = - \sum_i \log P(c_i | c_{i-j}, \dots, c_{i-1}; \Theta) \quad (1)$$

where  $j$  is the window length of previous tokens and  $\Theta$  means parameters of the LLM.

After being pre-trained, LLMs are adapted to a specific downstream task. At the beginning stage, LLMs are used in a fine-tuning manner, which are continually optimized on specific code generation datasets. With the size of LLMs growing rapidly such as ChatGPT [1] and CodeGen [47], fine-tuning is neither economical nor practical, in contrast, a convenient solution ICL arises.

## 2.2 In-Context Learning

In-context Learning (ICL) refers to an emerging ability of LLMs. Formally, given a set of requirement-code examples  $T = \{x_k, y_k\}_{k=1}^m$ , a test requirement  $x_t$  and a LLM with frozen parameters  $\Theta$ , ICL defines the generation of a program  $y_t$  as follows:

$$y_t \sim P(y_t | \underbrace{x_1, y_1, \dots, x_m, y_m}_{\text{context}}, x_t, \Theta) \quad (2)$$

where  $\sim$  represents decoding strategies such as greedy decoding and nuclear sampling [22] in code generation. The generation procedure is attractive as the parameters of LLMs are not need to be updated when executing a new task, which is efficient and practical.

As demonstrated in Formula 2, LLMs learn the task and generate programs depending on the selected examples. The performance of ICL can vary from almost random to near the state-of-the-art approach due to the different qualities of in-context examples. Thus, selecting appropriate examples is a significant research topic. In this paper, we propose a novel learning-based approach to select in-context examples for code generation.

## 3 METHOD: LAIL

In this section, we propose an effective learning-based approach LAIL to select examples for ICL in code generation, as shown in Figure 2. For a candidate example, LAIL exploits LLMs themselves to estimate it by incorporating the prediction probability of the ground-truth program given a requirement and the example, then labels examples as positive and negative according to their probability feedback (Section 3.1). Our approach imports a contrastive learning objective to optimize a neural retriever, resulting in alignment with the preference of LLMs (Section 3.2). In the inference stage, given a test requirement, LAIL selects a set of examples as a prompt and a LLM generates programs with the prompt (Section 3.3).

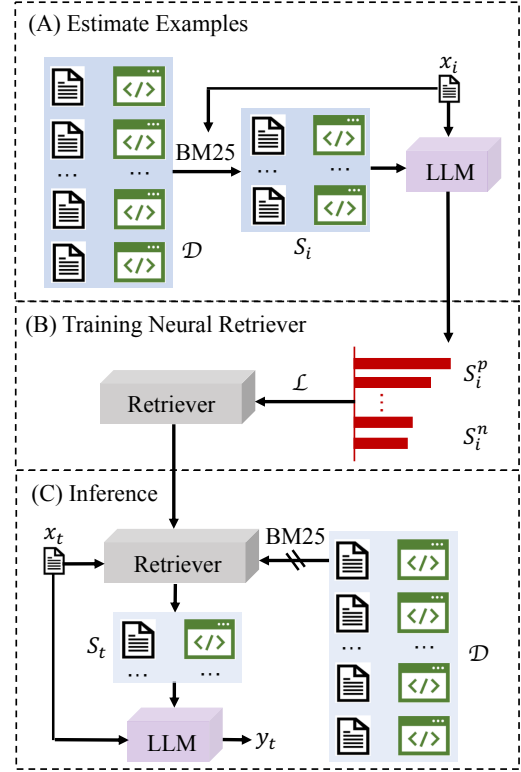
### 3.1 Estimate and Label Examples

Different LLMs are equipped with diverse prior knowledge, thus they have disparate preferences for candidate examples in code generation. To mitigate this phenomenon, we argue that a good selection approach should align with the preference of LLMs. Given a candidate example, we exploit LLMs themselves (LLM-aware) to measure it by incorporating the prediction probability of the ground-truth program given a requirement and the candidate item. Based on the probability feedback, we label candidate examples as positive and negative that are then used to train the retriever in Section 3.2.

Formally, given a requirement in the candidate pool, we use LLMs themselves to estimate whether each other candidate in the pool is beneficial to generate the ground truth of the requirement, and label each candidate based on its beneficial degree. Following previous studies [19, 31], this paper uses the training set as the candidate pool. The goal of this procedure can be formulated as follows:

$$\mathcal{D} = \{(e_i, \{S_i^p, S_i^n\})\}_{i=1}^N \quad (3)$$

where  $e_i = (x_i, y_i)$  is the  $i$ -th examples in the training set.  $S_i^p$  and  $S_i^n$  are the positive and negative example set of  $e_i$ , respectively.  $N$  is the number of examples in the training set  $\mathcal{R}$ .



**Figure 2: The overview of our LAIL. LAIL use LLMs themselves to estimate candidate examples and label them as positive and negative (A). Based on the label date, LAIL then trains a retriever to align with the preference of LLMs with a contrastive loss (B). Given a test requirement, the optimized retriever selects several examples as a prompt that is inputted to LLMs for code generation (C).**

Modeling on the full space of the training set is quadratic and thus prohibitive. To mitigate this limitation, we introduce a two-stage process to estimate examples in the training set. In the first stage, for each example  $e_i$ , we use an off-the-shelf retriever to select a set of examples  $S_i = \{(x_q^i, y_q^i)\}_{q=1}^t$  from  $\mathcal{R}$  where  $t \ll N$ . In this paper, we apply BM25<sup>1</sup> [42] to construct  $S_i$ , which measures textual n-gram matching between  $x_i$  and  $x_q^i$ . According to BM25 scores, the top- $r$  examples with higher BM25 scores constitute the set  $S_i$ . In the second stage, we leverage LLMs themselves to estimate each example in  $S_i$ . Concretely, we feed each example  $(x_q^i, y_q^i)$  in  $S_i$  and  $x_i$  into the LLM, and acquire the prediction probability of ground-truth program  $y_i$ . To quantify LLMs' probability feedback, we design a metric  $\mathcal{M}$  which is defined as follows:

$$\mathcal{M} = \frac{1}{|y_i|} \times P_{LLM}(y_i | x_q^i, y_q^i, x_i) \quad (4)$$

<sup>1</sup>We also use random and semantics-based approaches to filter examples and find that BM25 is more effective and efficient.

$$P_{LLM}(y_i|x_q^i, y_q^i, x_i) = \sum_{u=1}^{|y_i|} \log(p(t_{i,u}|x_q^i, y_q^i, x_i, t_{i,<u})) \quad (5)$$

where  $y_i = \{t_1, \dots, t_b\}$  and  $t_u$  is the  $u$ -th token in  $y_i$ . The metric can reflect how helpful the example  $(x_q^i, y_q^i)$  is for generating the target program  $y_i$ , which represents the preference of the LLM.

We then rank all examples in the set  $S_i$  according to their metric scores  $\mathcal{M}$  from high to low. The top- $z$  examples with higher scores are inputted into the positive example set  $S_i^p$ , meanwhile, the bottom- $v$  examples constitute the negative example set  $S_i^n$  since good examples should be beneficial for LLMs to generate correct programs. We apply this two-stage procedure to the entire training set and finally acquire the labeled data  $\mathcal{D}$ .

In our labeled data, the more helpful examples are labeled as positive examples and the examples with lower metric scores are treated as negative examples, which can reflect the preference of LLMs to candidate examples given a specific requirement.

### 3.2 Training Neural Retriever

As described in Section 3.1, the labeled data can reflect the preference of LLMs. In this section, we use the labeled data  $\mathcal{D}$  to train a neural retriever, aiming to align with the preference of LLMs. After being trained, given a test requirement, the retriever can select a set of candidate examples as a prompt that is beneficial for LLMs to generate correct programs.

To optimize the retriever, we import a contrastive learning objective. The objective targets to draw a test requirement with helpful training examples together and push the given requirement with knowledge-limited examples apart. Formally, for an example  $e_i$  in the labeled data  $\mathcal{D}$ , we randomly select an item  $e_i^p$  from  $S_i^p$  as its corresponding positive example. Meanwhile, we randomly pick a training example  $e_i^n$  from the set  $S_i^n$  and select  $h$  examples from the set  $\mathcal{H}$  (e.g.,  $\mathcal{H} = \mathcal{R} \setminus S_i$ ) as its negative example set  $\mathbb{N}_i$ . Next, We apply GraphCodeBERT [21] to encode their requirements and acquire corresponding representations. Then, we model the relations of these requirements with contrastive loss. Following SimCLR [15], the learning objective  $\mathcal{L}$  is formulated as:

$$\mathcal{L} = - \left[ \log \frac{e^{s(E_{[CLS]}^i, E_{[CLS]}^{i,p})/\tau}}{\sum_{E_{[CLS]}^{i,n} \in \mathbb{N}_i} e^{s(E_{[CLS]}^i, E_{[CLS]}^{i,n})/\tau}} \right] \quad (6)$$

where  $\tau$  is a temperature parameter.  $E_{[CLS]}$  means the entire representation of a requirement.  $s(\cdot)$  represents the cosine similarity of two vectors.

Based on the objective, our retriever learns the preference of LLMs, which will be effective to select helpful examples from the training set given a test requirement, resulting in promoting LLMs generating more correct programs.

### 3.3 Inference

During inference, instead of using heuristical approaches, we apply LAIL to select limited examples from the training set, which can retrieve examples with high metric scores given a test requirement and thus are helpful for LLMs in code generation. Specifically, we first feed requirements in the training set  $\mathcal{R}$  to our trained retriever

respectively and acquire their representational vectors  $\{E_{CLS}^i\}_{i=1}^N$ . Given a test requirement  $x_t$ , we obtain its representation  $E_{CLS}^t$  by the retriever. Next, we match  $E_{CLS}^t$  and  $E_{CLS}^i$  and thus lead to  $N$  pairs of representations  $\{(E_{CLS}^t, E_{CLS}^i)\}_{i=1}^N$ . We calculate their cosine similarities  $\{c_i\}_{i=1}^N$  of all pairs and rank candidate examples according to similarity scores from high to low. Then, the top- $r$  examples are selected for ICL. We concatenate the top- $r$  examples as a prompt  $\{e_1, \dots, e_i, \dots, e_r\}$  where their similarity scores gradually decrease (e.g.,  $c_1 < c_i < c_r$ ). We feed the prompt and the test requirement into LLMs and make LLMs generate programs with nuclear sampling [22] as described in Formula 2.

Note that LAIL only needs to encode training examples once. Given a test requirement, LAIL just calculates cosine similarities between it and all candidates. Thus, the efficiency of our approach is acceptable compared to other heuristic approaches.

## 4 STUDY DESIGN

To investigate the effectiveness of our LAIL, we perform a large-scale study to answer three research questions. In this section, we describe the details of our study, including datasets, evaluation metrics, baselines, base LLMs, and experimental details.

### 4.1 Research Questions

Our study aims to answer the following research questions.

**RQ1: How does LAIL perform compared to the state-of-the-art baselines?** This RQ aims to verify that LAIL can generate more correct programs than state-of-the-art (SOTA) baselines. We apply three LLMs to evaluate our approach. We compare LAIL to 7 baselines in three code generation datasets and employ unit tests to check the correctness of generated programs.

**RQ2: Do developers prefer programs generated by LAIL?** The ultimate goal of a code generation model is to assist developers in writing programs. In this RQ, we hire 10 developers to manually estimate the programs generated by our LAIL and baselines. We evaluate these programs in three aspects, including correctness, code quality, and maintainability.

**RQ3: What is the better design choice to estimate examples?** In section 3.1, we leverage the prediction probability of ground-truth programs to estimate candidate examples. In this RQ, we explore other design choices to measure examples by LLMs themselves and compare them to our design.

### 4.2 Datasets

We conduct experiments on three code generation datasets, including MBJP [8] in Java, MBPP [9] in Python, and MBCPP [8] in C++. The statistics of these datasets are given in Table 1.

**MBPP** [9] contains 974 Python programming problems constructed by crowd-sourcing. Each example consists of a brief description, a single self-contained function solving the problem specified, and three test cases to evaluate the correctness of the generated programs. The problems range from simple numeric manipulations or tasks that require the basic usage of standard library functions to tasks that demand nontrivial external knowledge.

**MBJP** [8] and **MBCPP** [8] are derived from MBPP [9]. MBJP and MBCPP contain 966 and 848 crowd-sourced programming problems

**Table 1: Statistics of three datasets on the different split sets.**

	M BJP	M BPP	M BCPP
Language	Java	Python	C++
#Train	383	384	413
#Dev	90	90	-
#Test	493	500	435
Avg. tests per example	3	3	3
Avg. tokens in requirement	16.71	16.50	17.38
Avg. tokens in code	247.79	92.68	113.94

in Java and C++, respectively. Each problem consists of a description, an individual function, and three test cases. These problems cover programming fundamentals, standard library functionality, etc.

We follow previous studies [8, 9] to split three datasets into the training set, the valid set, and the test set, respectively. We measure the performance of different ICL approaches on the test set.

### 4.3 Evaluation Metrics

Following prior code generation studies [8, 14], we leverage Pass@ $k$  to evaluate our approaches. Pass@ $k$  evaluates the functional correctness of the generated programs by executing test cases. Precisely, given a test requirement, we generate  $k$  programs using the sampling strategy. If any of the generated  $k$  programs passes all test cases, we think the requirement is solved. Finally, the percentage of solved requirements in all test requirements is treated as Pass@ $k$ . In this paper, we set  $k$  to 1, 3, and 5.

We notice that previous studies [20, 21] also use some match-based metrics such as BLEU [37] and CodeBLEU [41]. These metrics focus on the textual similarity or syntactic similarity between reference programs and generated programs. However, existing work [24] has shown that match-based metrics can not effectively measure the functionality of programs. In this paper, we follow recent studies [24, 29, 38] use execution-based metrics (e.g. Pass@ $k$ ).

### 4.4 Baselines

There are a few studies to investigate in-context example selection for code generation such as zero-shot learning [11], random selection [14], and AceCoder [31], where AceCoder [31] is the state-of-the-art (SOTA) baseline.

**Zero-Shot Learning** [11] directly inputs a requirement into LLMs without any examples as the prompt. LLMs generate programs for the given requirement.

**Random** [14] randomly select a few examples from the training set and construct a prompt. Then, LLMs predict the source code based on the prompt.

**AceCoder** [31] uses BM25 [42] to calculate the textual similarities between a test requirement and the requirements of candidates, and retrieve a set of examples with high similarities. In inference, it first generates test cases and then predicts programs. For a fair comparison, we use the same pipeline to LAIL and all baselines.

To extensively evaluate the effectiveness of our LAIL, we also transfer some advanced ICL approaches in natural language processing to the source code.

**TOP-k-SBERT** [35] leverages Sentence-BERT [40], a representative sentence encode, to encode all requirements in the training set. Given a test requirement, we first encode it and computer semantic similarities between it and the training requirements. Next, we select the top- $k$  similar examples.

**TOP-k-GraphCodeBERT** [21] is a variant of TOP-k-SBERT. It applies GraphCodeBERT to encode requirements and retrieve a few examples from the training set based on semantic similarity.

**TOP-k-VOTE** [43] is a graph-based method [43], to votes examples. It first encodes each example by GraphCodeBERT and each example is as a vertice in the graph. Each vertice connects with its  $k$  nearest vertices based on their semantic similarities. Finally, it treats the  $k$  nearest examples as a prompt.

**Uncertainty-Target** [17] assumes examples with higher uncertainty have a greater impact on LLMs. It defines uncertainty as the perplexity when LLMs generate ground truths. The approach computes the uncertainty of each candidate and selects  $k$  items with high perplexity as a prompt.

## 4.5 Base Large Language Models

This paper focuses on code generation with LLMs. Thus, we select three recently proposed LLMs for code generation as the base models. The details of the base models are shown as follows.

**ChatGPT** [1] is the state-of-the-art language model for code generation. It is trained on a large amount of natural language text and programming data. Then, ChatGPT is continually trained with reinforcement learning and learns to align with human instructions. We leverage OpenAI’s APIs to access ChatGPT (*i.e.*, gpt-3.5-turbo).

**GPT-3.5** [2] is a powerful large language model and is trained on large of unlabeled corpus. In this paper, we use OpenAI’s APIs to access the latest version with 175 billion parameters (*i.e.*, text-davinci-003).

**CodeGen** [47] is a family of language models for code generation. CodeGen is training with a 635GB code corpus and 1159GB English text data. In this paper, we leverage the largest version with 16 billion parameters (*i.e.* CodeGen-Multi-16B).

## 4.6 Implementation Details

**Estimate and Label Examples.** We use greedy decoding to generate programs and collect the predicted probabilities of ground-truth programs. Following the previous generation works [31], we set the max generated length to 500 tokens. The number of examples in the set  $S_i$  is 50 for efficiency. Note that the size of  $S_i$  larger the better performance of LLMs and we leave other sizes in our future work. We set  $z$  and  $v$  as 5 respectively and analyze their effects on performance in Section 6.3. Thus, the positive set  $S_i^p$  contains 5 examples with higher metric scores and the negative set  $S_i^n$  has 5 items with lower scores.

**Training Neural Retriever.** In this procedure, we set the number of the negative example set  $\bar{N}_i$  as 64. In other words,  $h$  is set to 63. For each epoch, we randomly select 63 examples from the set  $\mathcal{H}$ . We also attempt to set the size of  $\bar{N}_i$  to 32 and 128 in Section 6.3. The learning rate is 5e-5 and the batch size is 32. We train the retriever for about 1 hour on 4 NVIDIA A100.

**Inference.** We treat a large language model as a black-box generator and sample programs from it. The input of LLMs only contains

**Table 2: Evaluation results of our LAIL and baselines on CodeGen at the three code generation datasets. Numbers in bold indicate that the improvement is significant and the percentages in red color mean the relative improvements compared to the SOTA baseline.**

	MBJP			MBPP			MBCPP		
	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5
Zero-Shot Learning [11]	14.27	23.69	27.83	8.80	20.60	25.60	15.39	25.97	30.94
Random [14]	15.74	24.25	28.14	15.20	25.40	28.80	15.70	28.72	32.25
AceCoder [31]	17.65	27.18	30.63	16.80	26.40	29.40	17.47	30.11	33.78
TOP-k-SBERT [35]	16.63	26.77	29.21	15.40	25.00	29.00	18.09	29.43	33.33
TOP-k-GraphCodeBERT [21]	17.44	24.34	28.40	17.40	25.80	28.40	18.16	30.48	35.47
TOP-k-VOTE [43]	15.42	21.70	23.73	17.40	26.00	29.40	17.01	30.03	35.63
Uncertainty-Target [17]	17.09	23.06	27.93	14.60	24.20	28.80	17.70	30.12	34.62
LAIL	<b>21.30</b>	<b>28.49</b>	<b>32.05</b>	<b>18.60</b>	<b>27.80</b>	<b>30.60</b>	<b>19.08</b>	<b>31.36</b>	<b>37.94</b>
Relative Improvement	11.58%	4.82%	4.64%	6.89%	6.92%	4.08%	5.07%	2.89%	6.96%

a prompt and a test requirement without any natural language instructions. During the sampling, we use nuclear sampling [22] to decode, where the temperature is 0.8 and the top-p is 0.95. The maximum generated length is 500 tokens. For each test requirement, we generate 5 programs. For a fair comparison, we set the same parameters to generate programs for all baselines and our LAIL.

## 5 RESULTS AND ANALYSIS

### RQ1: How does LAIL perform compared to the state-of-the-art baselines?

In RQ1, we apply our LAIL and baselines to two LLMs and measure the correctness of generated programs.

**Setup.** We compare our LAIL and baselines (Section 4.4) on three code generation datasets (Section 4.2). The evaluation metric is Pass@k as described in Section 4.3. For the metric, higher scores indicate better performance of approaches.

**Results.** Table 2 and Table 3 report the Pass@k ( $k \in [1, 3, 5]$ ) of different approaches on CodeGen and GPT-3.5, respectively. Numbers in bold mean that the improvement is significant compared with baselines, and the percentages in red color represent the relative improvements compared to the state-of-the-art (SOTA) baseline.

**Analysis.** (1) LAIL achieves the best performance among all approaches with significant improvements. In all datasets, LAIL generates more correct programs than baselines. Compared to the SOTA baseline, LAIL outperforms it by 11.58%, 6.89%, and 5.07% on CodeGen, and acquires 4.38%, 2.85%, and 2.74% improvement on GPT-3.5 in Pass@1, respectively. Note that Pass@1 is a very strict metric and is difficult to be improved. The significant improvements prove the superiority of our LAIL in code generation. (2) Selecting proper examples is critical to the performance of ICL. Compared to random selection, LAIL acquires 8.31%, 30.60%, and 32.86% absolute improvements in Pass@1 on GPT-3.5. AceCoder and other heuristic approaches further improve code generation performance by selecting textually or semantically similar examples. LAIL exploits LLMs themselves to estimate and trains a neural retriever to align with the preference of LLMs, then achieves the best results among all approaches. That demonstrates the importance of examples in ICL and verifies the reasonability of using

LLMs themselves to label examples. (3) LAIL is effective in LLMs with different sizes and different programming languages. As above described, our approach on CodeGen and GPT-3.5 both achieve the best performance among all approaches. Table 2 and Table 3 also show that LAIL can generate more correct programs on all datasets including MBPP (Python), MBJP (Java), and MBCPP (C++). This reveals LAIL is generalized and can be applied to different LLMs and languages.

**Answer to RQ1: LAIL achieves the best results among all baselines. In three datasets, LAIL acquires 11.58%, 6.89%, and 5.07% improvements on CodeGen, and achieves 4.38%, 2.85%, and 2.74% improvements on GPT-3.5 at Pass@1. The significant improvements prove our approach can align with the preference of LLMs. Thus, it is able to select suitable examples for ICL and generates more correct programs.**

### RQ2: Do developers prefer programs generated by LAIL?

The goal of a code generation approach is to assist developers in writing programs. Thus, a good program not only satisfies the requirement but also is easy to read and maintain. In this question, we manually verify generated programs in three aspects.

**Setup.** Following previous work [29], we manually measure programs generated by different approaches in three aspects (e.g. correctness, code quality, maintainability). Correctness measures whether a program satisfies the given requirement, code quality verifies whether a program does not contain bad code smell, and maintainability measures whether the implementation is standardized and easy to read. For each aspect, the score is an integer and ranges from 0 to 2. The higher score, the better code is.

Specifically, we randomly select 50 test samples from MBPP, MBJP, and MBCPP, respectively<sup>2</sup>. Then, we use LAIL and baselines to generate programs for these examples. Finally, we obtain 400 ( $50 \times 8$ ) programs for human evaluation. The 400 programs are randomly divided into 5 groups. We hire 10 developers to verify

<sup>2</sup>The confidence level is 85% and the margin of error is 15% with 50 selected examples.

**Table 3: The performance of our LAIL and the existing approaches on GPT-3.5 at the three code generation datasets. The numbers in red color represent the relative improvements compared to the SOTA baseline.**

	MBJP			MBPP			MBCPP		
	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5
Zero-Shot Learning [11]	44.83	53.05	59.72	20.00	27.00	29.60	21.85	37.94	49.90
Random [14]	47.87	59.83	63.69	43.00	56.40	60.80	50.02	61.43	65.28
AceCoder [31]	50.91	61.25	65.15	47.40	61.00	64.20	53.25	63.22	66.21
TOP-k-SBERT [35]	50.30	60.46	64.70	48.00	61.40	64.00	52.87	63.21	67.13
TOP-k-GraphCodeBERT [21]	50.30	60.85	64.50	49.20	58.20	64.20	52.64	63.19	67.28
TOP-k-VOTE [43]	50.52	60.45	63.49	47.80	59.20	63.40	51.03	62.98	65.51
Uncertainty-Target [17]	47.67	59.23	63.69	37.80	51.40	57.80	42.87	52.60	57.18
LAIL	<b>53.14</b>	<b>62.87</b>	<b>66.72</b>	<b>50.60</b>	<b>62.40</b>	<b>65.20</b>	<b>54.71</b>	<b>65.98</b>	<b>69.67</b>
Relative Improvement	<b>4.38%</b>	<b>2.64%</b>	<b>2.41%</b>	<b>2.85%</b>	<b>1.63%</b>	<b>1.53%</b>	<b>2.74%</b>	<b>4.37%</b>	<b>3.55%</b>

**Table 4: The results of human evaluation. The numbers in red represent LAIL’ relative improvements compared to the SOTA baseline. “CR”, “QL”, and “MT” mean correctness, code quality, and maintainability, respectively.**

Approach	CR	QL	MT
Zero-Shot Learning [11]	0.472	1.079	1.183
Random [14]	0.925	1.416	1.520
AceCoder [31]	1.523	1.652	1.647
TOP-k-SBERT [35]	1.569	1.640	1.665
TOP-k-GraphCodeBERT [21]	1.582	1.738	1.609
TOP-k-VOTE [43]	1.325	1.624	1.631
Uncertainty-Target [17]	0.736	1.297	1.324
LAIL	<b>1.764</b>	<b>1.839</b>	<b>1.782</b>
Relative Improvement	<b>11.504%</b>	<b>5.811%</b>	<b>7.027%</b>

these programs. These developers are computer science students or industrial practitioners and write programs at least for 3 years. Each group is measured by two developers and the final score is the average of two developers’ scores.

**Results.** The results of the human evaluation are shown in Table 4. The percentages in parentheses are the relative improvements compared to the SOTA baseline.

**Analyses.** (1) Developers prefer programs generated by LAIL over all baselines. LAIL substantially outperforms all baselines in three aspects, which demonstrates the effectiveness of our approach. (2) In addition to satisfying requirements, programs provided by our approach are easier to read and have less bad smell. Particularly, our LAIL surpasses the SOTA approach by 5.811% in code quality and 7.027% in maintainability. (3) Some heuristic methods that select the same examples as prompts for all test requirements are not optimal for ICL. As shown in Figure 4, Uncertainty-Target and Vote-K perform comparably to the random approach, which indicates that applying the same prompt to all test data is not a good choice, and different test requirements should use their appropriate prompts.

**Answer to RQ2: Human evaluation shows that LAIL significantly outperforms SOTA baselines on all three aspects. Programs generated by LAIL can better satisfy requirements, be easier to read, and have little bad code smell. That is, developers prefer programs generated by LAIL among all approaches.**

### RQ3: What is the better design choice to estimate examples?

In this paper, we apply the predicted probability of ground truths to estimate examples. We further investigate two plausible choices to measure candidates by LLMs themselves.

**Setup.** We design two plausible approaches including the Match-BLEU and the Match-CodeBLEU formats. The Match-BLEU format uses the BLEU score of the generated program to measure examples. The Match-CodeBLEU approach applies the CodeBLEU score of the predicted program to label examples. In this RQ, we select a representative LLM (*i.e.* GPT-3.5) as the base model and evaluate their performance on three datasets.

**Results.** The results of different feedback scores are represented in Table 5. For comparison, we also present the results of the existing SOTA baseline in Table 5.

**Analyses.** (1) Our probability-based approach is better than Match-BLEU and Match-CodeBLEU. In all datasets, our LAIL outperforms them by 4.05%, 5.86%, and 3.77% on three datasets at Pass@1, respectively. We argue that the predicted probability of ground truth accurately reflects the certainty of LLMs for correct programs. The BLEU-based and CodeBLEU-based methods only reflect the literal accuracy of LLMs’ prediction, but can not provide how certain LLMs are when predicting programs. (2) Match-CodeBLEU approach is more effective than Match-BLEU method. Match-CodeBLEU outperforms Match-BLEU on all datasets. The reason might be that CodeBLEU measures the n-gram match, the semantic match, and the syntactic match between the hypothesis programs and the reference code tokens, which is better to evaluate the generated programs than BLEU.

**Answer to RQ3: Probability-based approach is better than Match-CodeBLEU and Match-BLEU formats. It can effectively reflect the LLMs' preference for candidate examples given a test requirement. .**

## 6 DISCUSSION

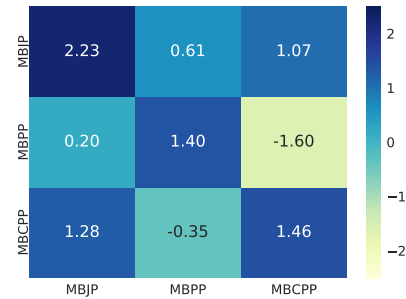
### 6.1 Transferability

We explore whether our retriever based on one LLM's feedback and specific dataset can be transferred to other LLMs or code generation datasets without further tuning. This is a significant research question since the retriever for each LLM and dataset needs to be trained in real applications.

**6.1.1 Transfer across LLMs.** We consider transferring the retriever based on one LLM's feedback to another LLM. Specifically, we use a source LLM (e.g. CodeGen or GPT-3.5) to estimate examples for training a retriever and then apply the retriever to another target LLM (e.g. ChatGPT) in generating programs. Table 6 shows the performance of ChatGPT in three datasets. We surprisingly find our retriever based on CodeGen and GPT-3.5 can bring obvious improvements to ChatGPT. In particular, in terms of Pass@1, ChatGPT achieves 2.56% improvements from CodeGen's feedback and 4.31% enhancements from GPT-3.5's feedback compared to the SOTA baseline. The phenomena demonstrate that our approach has satisfying transfer ability across different LLMs. Note that ChatGPT can not provide the prediction probability of ground truths in practice, thus LAIL is a quite meaningful approach, especially for LLMs whose parameters are unavailable. Besides, the performance of ChatGPT from GPT-3.5's feedback is higher than the counterpart from CodeGen's feedback. The reason might be that GPT-3.5 and ChatGPT have comparable abilities in code generation, thus their preference is similar and GPT-3.5 can provide more proper examples as prompts.

To verify the transfer ability among LLMs with different sizes, we further evaluate the performance of CodeGen based on GPT-3.5's feedback and the results of GPT-3.5 from CodeGen's feedback, where the parameter size of CodeGen is much smaller than the counterpart of GPT-3.5. As shown in Table 6, compared to the random approach, the retriever learned from CodeGen achieves 8.25%, 12.09%, and 4.32% relative improvements to GPT-3.5 on Pass@1, meanwhile our retriever under GPT-3.5's feedback brings 22.30%, 14.47%, and 16.94% improvements to CodeGen, respectively. This indicates our approach has satisfying transfer ability among LLMs and can bring improvements to target LLMs.

**6.1.2 Transfer across datasets.** Considering that the compositional features of natural language are general, the retriever trained on one dataset may apply to other datasets and exploit similar knowledge in different datasets. In this section, we further investigate whether a retriever trained on one dataset can transfer to others. We transfer the retriever among three datasets (e.g., MBJP, MBPP, and MBCPP) and Figure 3 demonstrates the transferring results on Pass@1. We find that most retrievers can successfully transfer to other datasets and bring improvements compared to their SOTA baselines. Concretely, the retriever trained on MBJP (MBCPP)



**Figure 3: Results of transferring the retriever trained on one dataset (row) to others (column) on GPT-3.5 in MBPP dataset.**

achieves 1.28% (1.07%) absolute improvements when it migrates to MBCPP (MBJP). Meanwhile, the MBCPP-trained retriever hardly transfers to MBPP and the MBPP-based retriever suffers the generation performance on MBPP. The reason might be that Java and C++ are object-oriented programming languages, and their syntax and code morphology are similar. Exploring a retriever that is suitable for many datasets is a challenging but meaningful research question, and we leave this topic as our future work.

### 6.2 Impacts of In-Context Example Numbers

Most of the LLMs are trained with a limited input length, which restricts the number of examples in the prompt. Previous studies [19, 35] also find that LLMs are affected by the number of in-context examples. Here we explore the impacts of the number of examples in the prompt on baselines and our approach. Figure 4 reports how the performance of LAIL and GraphCodeBET change with respect to different in-context example numbers in MBPP dataset. Note that GraphCodeBET is the best baseline in this dataset, thus we choose it to compare with our approach. We observe that the performances of GraphCodeBET and our approach monotonically increase in rough with the increase of in-context example numbers. In addition, our LAIL always outperforms GraphCodeBERT in all cases, which further proves the superiority of our approach even with different in-context example numbers.

### 6.3 Effects of Contrastive Learning Parameters

The number of negative examples  $\mathbb{N}_i$  (named  $\mathbb{N}_i$  for convenience in Figure 7) is an important element in contrastive learning for training a neural retriever. We investigate how the element affects the performance of our approach. As shown in Figure 7, with the increasing of negative example numbers, LAIL achieves consistent improvements and can generate more correct programs. Besides, as described in 3.2, we randomly select an example  $e_i^n$  from  $S_i^n$  into  $\mathbb{N}_i$ . We also explore the effect of the number of selected examples from  $S_i^n$  (dubbed  $\tau_{ne}$ ). We can find that the more hard negative examples, the worse our approach performs. We argue that the examples of  $S_i^n$  are with high BM25 scores among all candidate examples and thus  $S_i^n$  may contain some false negative examples. Selecting more examples from  $S_i^n$  will affect a retriever learning for choosing suitable in-context examples. Between the two factors,



**Table 5: The comparison of different formats to estimate examples on three datasets. The values in parentheses mean significant improvements. "Probability-Based" represents our LAIL.**

	MBJP			MBPP			MBCPP		
	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5
Random [14]	47.87	59.83	63.69	43.00	56.40	60.80	50.04	61.45	65.28
Match-BLEU	50.19	59.87	63.92	46.20	60.20	62.40	50.26	61.47	65.73
Match-CodeBLEU	51.07	60.04	64.58	47.80	59.20	63.80	52.72	62.03	66.80
Probability-Based (LAIL)	<b>53.14</b>	<b>62.87</b>	<b>66.72</b>	<b>50.60</b>	<b>62.40</b>	<b>65.20</b>	<b>54.71</b>	<b>65.98</b>	<b>69.67</b>
Relative Improvement	4.05%	4.71%	3.31%	5.86%	3.65%	2.19%	3.77%	6.37%	4.30%

**Table 6: Results of transferring a retriever learned on one LLM to others on three datasets. ♣ means the source LLM.**

♣ ChaTGPT	MBJP			MBPP			MBCPP		
	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5	Pass@1	Pass@3	Pass@5
Zero-Shot Learning [11]	16.63	34.48	44.21	26.60	32.00	34.60	30.34	57.01	63.68
Random [14]	53.34	62.27	65.72	50.80	60.60	63.40	40.15	60.06	66.28
AceCoder [31]	54.46	64.01	66.75	54.40	62.40	65.20	42.53	62.29	68.32
TOP-k-SBERT [35]	54.26	63.89	66.53	54.80	63.00	65.40	43.37	61.45	67.11
TOP-k-GraphCodeBERT [21]	53.95	62.67	66.32	54.20	62.60	65.20	44.08	61.61	68.27
TOP-k-VOTE [43]	51.93	62.88	65.72	42.00	56.00	61.00	42.74	62.68	67.73
Uncertainty-Target [17]	49.69	60.45	54.50	36.40	51.80	56.80	42.46	61.84	66.67
CodeGen									
LAIL	<b>54.79</b>	<b>64.70</b>	<b>67.43</b>	<b>55.60</b>	<b>63.80</b>	<b>66.20</b>	<b>45.21</b>	<b>63.45</b>	<b>68.93</b>
Relative Improvement	0.61%	1.08%	1.02%	1.44%	1.27%	1.22%	2.56%	1.23%	0.89%
GPT-3.5									
LAIL	<b>55.97</b>	<b>64.97</b>	<b>68.27</b>	<b>56.20</b>	<b>64.60</b>	<b>66.80</b>	<b>45.98</b>	<b>63.84</b>	<b>70.35</b>
Relative Improvement	2.77%	1.50%	2.28%	2.56%	2.54%	2.14%	4.31%	1.85%	2.97%
♣ GPT-3.5									
CodeGen									
Random [14]	47.87	59.83	63.69	43.00	56.40	60.80	50.02	61.43	65.28
LAIL	<b>51.82</b>	<b>61.75</b>	<b>64.89</b>	<b>48.20</b>	<b>59.00</b>	<b>64.20</b>	<b>52.18</b>	<b>64.54</b>	<b>67.90</b>
Relative Improvement	8.25%	3.21%	1.88%	12.09%	4.61%	5.59%	4.32%	5.06%	4.01%
♣ CodeGen									
GPT-3.5									
Random [14]	15.74	24.25	28.14	15.20	25.40	28.80	15.70	28.72	32.25
LAIL	<b>19.25</b>	<b>27.53</b>	<b>30.98</b>	<b>17.40</b>	<b>26.00</b>	<b>30.00</b>	<b>18.36</b>	<b>30.41</b>	<b>35.54</b>
Relative Improvement	22.30%	13.53%	10.09%	14.47%	2.36%	4.17%	16.94%	5.88%	10.20%

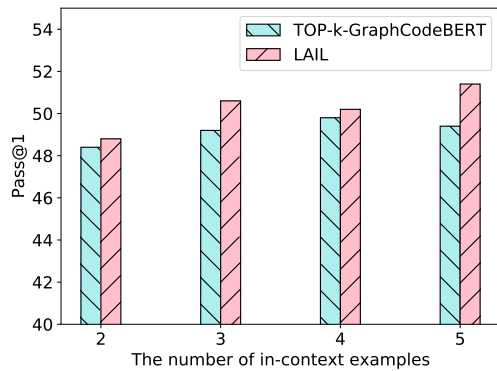
the number of selected examples from  $S_i^n$  has a greater influence on the performance of our approach.

## 6.4 Threats to Validity

There are two main threats to the validity of our work.

**The generalizability of our experimental results.** For the datasets, we follow previous studies [8, 9, 14] and leverage three representative code generation datasets. The three datasets cover different programming languages (e.g., Java, Python, and C++) and

come from real-world software communities. To verify the superiority of LAIL, we consider seven existing ICL approaches in both the code generation task and many maintain natural language tasks. In addition, to effectively evaluate our approach, we select a series of advanced pre-trained large language models (CodeGen [47], GPT-3.5, and ChatGPT) as base models in the past three years. We apply our approach and baselines to base models and evaluate their performance in code generation. For the metric, following existing works [13, 47], we select a widely used Pass@k metric to evaluate



**Figure 4: The performance of the different number of in-context examples on GPT-3.5 in MBPP datasets.**

**Table 7: Effects of parameters.**

	Pass@1	Pass@3	Pass@5
Negative examples ( $N_i$ )			
32	50.20	61.80	65.00
64	50.60	62.40	65.20
128	50.80	62.40	65.40
$\tau_{ne}$			
1	50.60	62.40	65.20
5	50.20	62.20	65.00
10	49.40	61.80	64.40

all approaches. It is an execution-based metric that utilizes test cases to check the correctness of generated programs. To ensure fairness, we execute each method three times and report the average experimental results.

**The implementation of models and prompts.** It is widely known that deep neural models are sensitive to the implementation details. In this paper, we need to execute all baselines and our approach on three base models. For baselines, we apply the source code and parameters published by their original papers [21, 40, 47]. For LLMs (e.g., GPT-3.5 and ChatGPT), the hyper-parameters (e.g., temperature) of sampling will impact their outputs. In experiments, we follow previous studies [29, 47] to set hyper-parameters for all approaches. In addition, the performance of LLMs heavily depends on the prompts, including the instruction and the number of examples. To alleviate this threat, we leverage the same number of examples for all approaches and directly construct prompts without any natural language instructions. Besides, it’s worth noting that the candidate pool for example selection also influences ICL performance. A large-scale study on 13.2 million real code files showed the proportion of reused code is up to 80% [36]. Thus, we use the training set of each dataset as the candidate set for all approaches in this paper and believe that the training set can provide supporting examples for test requirements. We do not tune the prompt and hyper-parameters experimentally and set them empirically. Thus, there might be room to tune hyper-parameter settings of our approach for more improvements.

## 7 RELATED WORK

### 7.1 Pre-trained Language Models

Code generation aims to automatically generate programs given natural language requirements. Pre-trained language models have achieved state-of-the-art results in code generation. They are pre-trained on large-scale unlabeled code corpus and then transferred to code generation. Existing pre-trained models can be divided into encoder-decoder models and decoder-only models.

**Encoder-decoder models** consist of an encoder and a decoder. An encoder takes a requirement as the input, and a decoder outputs a program. Many popular encoder-decoder architectures in natural language processing have been applied to source code, which often use denoising-based pre-training objectives to pre-train, such as masked sequence prediction, identifier tagging, and masked identifier prediction. For example, CodeT5 [44] is the variant of T5 model [39] to support programs. [5] is adapted from the BART model [28], which is pre-trained on a number of programs. Besides, some studies [12] further consider the naturalized feature of code and train models to generate natural programs based on noised programs.

**Decoder-only models** contain decoder networks, which are pre-trained with the next token prediction objective. Inspired by the success of GPT series [38] in natural language processing, researchers attempt to adapt similar models to source code. CodeGPT [34] is pre-trained on CodeSearchNet [23] corpus with the same setting and structure as GPT-2. CodeX is continually fine-tuned on GPT-3 [11] on code corpus from GitHub. CodeX is proficient in over a dozen languages (e.g., JavaScript and Python) and supports a commercial application (e.g., Copilot). Since its parameters are not available, many researchers try to replicate them and bring CodeParrot [3], GPT-CC [4], PyCodeGPT [45] and CodeGen [47]. CodeGeeX [46] is pre-trained on a large code corpus and can generate more than 20 programming languages. Lately, ChatGPT is proposed by OpenAI and achieves impressive performance in code generation. In this paper, we select three representative networks including CodeGen [47], GPT-3.5 [2], and ChatGPT [1] as our base models.

### 7.2 In-Context Learning

In-context learning (ICL) is an emerging approach to using large language models (LLMs). By providing limited examples as a prompt, ICL empowers LLMs to learn a specific downstream task. ICL [11] is first proposed in natural language processing and achieves impressive performance in many tasks [19, 27] such as text generation and sentiment classification.

Inspired by the success of ICL in natural language processing, researchers attempt to adapt ICL to source code [10, 14, 16]. They design task-specific instructions with a set of examples to prompt LLMs and improve the performance on many tasks (e.g., code generation [16] and code repair [25]). The popularity of ICL also introduces instability in performance: given different sets of examples as prompts, LLMs’ performance usually varies from random to near state-of-the-art [33]. Some studies [16, 31] make efforts to mitigate this issue. Researchers [16] randomly select limited examples for ICL and verify the results in code generation. AceCoder [31] uses BM25 to select n-gram matching examples and further generates more correct programs. However, these approaches are based on

lexical features, which require developers to design heuristics and lead to sub-optimal performance. Compared to the existing methods, LAIL is a learning-based approach for selecting in-context examples. Instead of considering textual similarity, it applies LLMs themselves to measure candidate examples and uses contrastive learning to learn LLMs' preferences in code generation.

## 8 CONCLUSION

Due to the instability of ICL, there is an increasing demand for in-context example selection in code generation. This paper proposes a novel learning-based ICL approach LAIL. We exploit LLMs themselves to estimate examples by considering the generation probabilities of ground-truth programs given a requirement, and label these examples based on LLMs' feedback. We then introduce a contrastive learning objective to train a retriever, resulting in acquiring the preferences of LLMs. Experimental results on three code generation datasets demonstrate that our proposed approach achieves excellent performance. LAIL also shows satisfactory transferability across LLMs and datasets, showing that it is an efficient and practical approach to code generation. In the future, we will apply our approach to other LLMs and datasets.

## REFERENCES

- [1] Available. <https://openai.com/chatgpt>.
- [2] Available. <https://platform.openai.com/docs/models/gpt-3-5>.
- [3] Available. <https://huggingface.co/codeparrot/codeparrot>.
- [4] Available. <https://github.com/CodedotAI/gpt-code-clippy>.
- [5] WU Ahmad, S Chakraborty, B Ray, and KW Chang. 2021. Unified pre-training for program understanding and generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [6] Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. 2018. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 1–37.
- [7] Miltiadis Allamanis, Earl T Barr, Premkumar T. Devanbu, and Charles Sutton. 2018. A Survey of Machine Learning for Big Code and Naturalness. *ACM Comput. Surv.* 51, 4 (2018), 81:1–81:37. <https://doi.org/10.1145/3212695>
- [8] Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, et al. 2022. Multi-lingual Evaluation of Code Generation Models. *arXiv preprint arXiv:2210.14868* (2022).
- [9] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732* (2021).
- [10] Patrick Bareiß, Beatriz Souza, Marcelo d'Amorim, and Michael Pradel. 2022. Code generation tools (almost) for free? a study of few-shot, pre-trained language models on code. *arXiv preprint arXiv:2206.01335* (2022).
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [12] Saikat Chakraborty, Toufique Ahmed, Yangruibo Ding, Premkumar T Devanbu, and Baishakhi Ray. 2022. Natgen: generative pre-training by “naturalizing” source code. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 18–30.
- [13] Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. *arXiv preprint arXiv:2207.10397* (2022).
- [14] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [15] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [16] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128* (2023).
- [17] Cody Coleman, Christopher Yeh, Stephen Mussmann, Baharan Mirzasoleiman, Peter Bailis, Percy Liang, Jure Leskovec, and Matei Zaharia. 2019. Selection via proxy: Efficient data selection for deep learning. *arXiv preprint arXiv:1906.11829* (2019).
- [18] Li Dong and Mirella Lapata. 2016. Language to Logical Form with Neural Attention. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 33–43.
- [19] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234* (2022).
- [20] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiao Cheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).
- [21] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366* (2020).
- [22] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751* (2019).
- [23] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436* (2019).
- [24] Jeevana Priya Inala, Chenglong Wang, Mei Yang, Andres Codas, Mark Encarnación, Shuvendu Lahiri, Madanlal Musuvathi, and Jianfeng Gao. 2022. Fault-aware neural code rankers. *Advances in Neural Information Processing Systems* 35 (2022), 13419–13432.
- [25] Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan. 2023. Impact of code language models on automated program repair. *arXiv preprint arXiv:2302.05020* (2023).
- [26] Xue Jiang, Yihong Dong, Lecheng Wang, Qiwei Shang, and Ge Li. 2023. Self-planning code generation with large language model. *arXiv preprint arXiv:2303.06689* (2023).
- [27] Itay Levy, Ben Bogin, and Jonathan Berant. 2022. Diverse demonstrations improve in-context compositional generalization. *arXiv preprint arXiv:2212.06800* (2022).
- [28] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* (2019).
- [29] Jia Li, Yongmin Li, Ge Li, and Zhi Jin. 2023. Structured Chain-of-Thought Prompting for Code Generation. *CoRR* (2023). <https://arxiv.org/abs/2305.06599>
- [30] Jia Li, Yongmin Li, Ge Li, Zhi Jin, Yiyang Hao, and Xing Hu. 2023. SkCoder: A Sketch-based Approach for Automatic Code Generation. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*. IEEE, 2124–2135. <https://doi.org/10.1109/ICSE48619.2023.00179>
- [31] Jia Li, Yunfei Zhao, Yongmin Li, Ge Li, and Zhi Jin. 2023. AceCoder: Utilizing Existing Code to Enhance Code Generation. *CoRR* (2023). <https://arxiv.org/abs/2303.17780v3>
- [32] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science* 378, 6624 (2022), 1092–1097.
- [33] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What Makes Good In-Context Examples for GPT-3? *arXiv preprint arXiv:2101.06804* (2021).
- [34] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664* (2021).
- [35] 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. *arXiv preprint arXiv:2305.14128* (2023).
- [36] Audris Mockus. 2007. Large-scale code reuse in open source software. In *First International Workshop on Emerging Trends in FLOSS Research and Development (FLOSS'07: ICSE Workshops 2007)*. IEEE, 7–7.
- [37] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318.
- [38] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).
- [39] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21, 140 (2020), 1–67.
- [40] Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084* (2019).
- [41] Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. 2020. Codebleu: a method for

- automatic evaluation of code synthesis. *arXiv preprint arXiv:2009.10297* (2020).
- [42] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [43] Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. 2022. Selective annotation makes language models better few-shot learners. *arXiv preprint arXiv:2209.01975* (2022).
- [44] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859* (2021).
- [45] Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen, and Jian-Guang Lou. 2022. CERT: Continual Pre-Training on Sketches for Library-Oriented Code Generation. *arXiv preprint arXiv:2206.06888* (2022).
- [46] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568* (2023).
- [47] Maosheng Zhong, Gen Liu, Hongwei Li, Jiangling Kuang, Jinshan Zeng, and Mingwen Wang. 2022. CodeGen-Test: An Automatic Code Generation Model Integrating Program Test Information. *arXiv preprint arXiv:2202.07612* (2022).