# Instructive Code Retriever: Learn from Large Language Model's Feedback for Code Intelligence Tasks

### Jiawei Lu*
The State Key Laboratory of
Blockchain and Data Security,
Zhejiang University
Hangzhou, China
jiaweilu@zju.edu.cn

### Haoye Wang*
Hangzhou City University
Hangzhou, China
wanghaoye@hzcu.edu.cn

### Zhongxin Liu†‡
The State Key Laboratory of
Blockchain and Data Security,
Zhejiang University
Hangzhou, China
liu_zx@zju.edu.cn

### Keyu Liang
The State Key Laboratory of
Blockchain and Data Security,
Zhejiang University
Hangzhou, China
liangkeyu@zju.edu.cn

### Lingfeng Bao
The State Key Laboratory of
Blockchain and Data Security,
Zhejiang University
Hangzhou, China
lingfengbao@zju.edu.cn

### Xiaohu Yang
The State Key Laboratory of
Blockchain and Data Security,
Zhejiang University
Hangzhou, China
yangxh@zju.edu.cn

## ABSTRACT

Recent studies proposed to leverage large language models (LLMs) with In-Context Learning (ICL) to handle code intelligence tasks without fine-tuning. ICL employs task instructions and a set of examples as demonstrations to guide the model in generating accurate answers without updating its parameters. While ICL has proven effective for code intelligence tasks, its performance heavily relies on the selected examples. Previous work has achieved some success in using BM25 to retrieve examples for code intelligence tasks. However, existing approaches lack the ability to understand the semantic and structural information of queries, resulting in less helpful demonstrations. Moreover, they do not adapt well to the complex and dynamic nature of user queries in diverse domains. In this paper, we introduce a novel approach named Instructive Code Retriever (ICR), which is designed to retrieve examples that enhance model inference across various code intelligence tasks and datasets. We enable ICR to learn the semantic and structural information of the corpus by a tree-based loss function. To better understand the correlation between queries and examples, we incorporate the feedback from LLMs to guide the training of the retriever. Experimental results demonstrate that our retriever significantly outperforms state-of-the-art approaches. We evaluate our model's effectiveness on various tasks, i.e., code summarization, program synthesis, and bug fixing. Compared to previous state-of-the-art algorithms, our method achieved improvements of 50.0% and 90.0% in terms of BLEU-4 for two code summarization datasets, 74.6% CodeBLEU on program synthesis dataset, and increases of 3.6 and 3.2 BLEU-4 on two bug fixing datasets.

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence**; • **Software and its engineering**;

## KEYWORDS

Software Engineering, Large Language Models, In-Context Learning

## 1 INTRODUCTION

Code intelligence tasks aim to automate the analysis, understanding, and generation of code. Representative code intelligence tasks include code summarization, bug fixing, and program synthesis. These tasks are highly valued because of their potential to mitigate the time-consuming nature of manual efforts, such as annotation composition, debugging, and code creation [29, 34, 52]. As the complexity of software systems continues to increase, there is a growing need for efficient and automated approaches to understanding, modifying, and improving source code [17, 25, 44, 51].

In recent years, pre-trained models like CodeBERT [14], Graph-CodeBERT [19] and CodeT5 [53] have made significant strides across various code intelligence tasks. They are pre-trained on large-scale code corpora with diverse objects. However, when dealing with various code intelligence tasks, these pre-trained models need to be fine-tuned separately on datasets for each task, limiting their generalizability. Additionally, the parameter sizes of these pre-trained models are relatively, restricting the knowledge they

---
*Equal Contribution.
†Corresponding Author.
‡Also with Hangzhou High-Tech Zone (Binjiang) Institute of Blockchain and Data Security.

can capture [54] and thereby preventing them from achieving satisfactory performance across multiple tasks simultaneously [27].

Recently, large language models (LLMs) have attracted a lot of attention due to their impressive effectiveness on diverse tasks [21, 55, 58]. In contrast to traditional pre-trained models, LLMs increase parameter counts tenfold or even hundredfold [8] and have been trained on massive datasets of both natural language and programming languages [47]. Given appropriate prompts, LLMs demonstrate exceptional proficiency across diverse code intelligence tasks, obviating the necessity for parameter fine-tuning. However, prior research indicates that LLMs exhibit sensitivity to input prompts [41]. To better guide LLMs in performing code intelligence tasks, several studies have explored In-Context Learning (ICL) [3, 16, 38], which provides LLMs with several input-output pairs to showcase how to follow the instruction. As illustrated in Figure 1, the retriever retrieves a series of demonstrations from the training set and adds them before the query, thereby prompting the LLM to output the correct answer. Instead of fine-tuning the model for a task, ICL can achieve efficient knowledge transfer with just a few examples, making it effective for handling diverse tasks and adapting to rapidly changing environments.

Previous researches show that the selection of demonstrations is crucial to the performance of In-Context Learning (ICL) [16, 31]. Gao et al. [16] find that BM25 [45] can effectively retrieve ICL examples for various code intelligence tasks. Ahmed et al. [3] improve their prompts with designed semantic facts for code summarization after BM25 retrieved examples. They name their approach as ASAP.

However, BM25-based retrievers have limited capabilities for understanding rich structural information in programming languages, which are useful for code intelligence tasks [60]. Additionally, they lack the essential domain knowledge and semantic understanding ability to comprehend the correlation between queries and examples. These capabilities are important for code intelligence tasks because programs in different forms may have similar functionality. Moreover, ASAP is only applicable to specific datasets and tasks.

In this research, we present a novel approach named **I**nstructive **C**ode **R**etriever (ICR) to address the aforementioned limitations. We use a combined encoding of text and syntax trees to calculate the similarity of examples. Then we propose a tree-based loss function to train the retriever, allowing the retriever to consider both semantic and structural information. Additionally, we utilize feedback from the LLM to rank the quality of examples, enabling the retriever to leverage the LLM's domain knowledge and semantic understanding ability through contrastive learning.

Specifically, in the training phase, ICR initially utilizes syntax trees to capture structural information from the queries and examples. Then we leverage the inherent capabilities of LLMs to discern the usefulness of examples for each query, thereby obtaining both positive and negative instances for training a retriever capable of selecting superior examples. Finally, we train the retriever iteratively by minimizing a customized tree-based loss between the retriever's similarity and the sequence feedback from the large model. In the testing phase, we utilize the trained retriever to retrieve valuable examples. These instances are then used to construct prompts, which are employed to enhance the performance of LLMs with respect to the given code intelligence tasks.
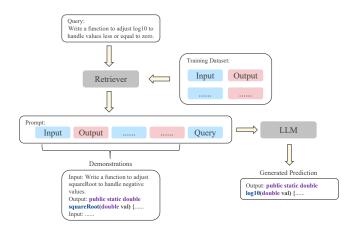


**Figure 1: An Example of In-Context Learning on Program Synthesis Task.**

We evaluate our proposed retriever ICR on three code intelligence tasks, i.e., code summarization, program synthesis, and bug fixing, using eight datasets in total. Compared to previous state-of-the-art algorithms, ICR achieves improvements of 50.0% and 90.0% on the Java and Python datasets in terms of BLEU-4 for code summarization, 74.6% in terms of CodeBLEU on program synthesis, and increases of 3.6 BLEU-4 and 3.2 BLEU-4 on two datasets for bug fixing with the base model GPT-Neo-2.7B. When transferring ICR trained on GPT-Neo-2.7B to Code Llama-13B, ICR improves over the state-of-the-art baselines by up to 18.3%, 20.5%, and 1.6% on the code summarization, program synthesis, and bug fixing tasks on BLEU-4, respectively. Additionally, ICR outperformed state-of-the-art methods on code summarization datasets in three languages that ICR had not been trained on. Moreover, we found that when we use high-quality examples, the order of examples has little impact on performance.

To summarize, the contributions of our work are:

- We propose a novel retriever to retrieve valuable demonstrations that can help large language models accomplish code intelligence tasks. To the best of our knowledge, this is the first method that utilizes large language model feedback for contrastive learning to enhance the performance of ICL in code intelligence tasks.
- To utilize the structural information from the corpus, we introduce a tree-based model loss function tailored specifically for code intelligence tasks.
- We evaluated the effectiveness of our approach across multiple programming languages and various code intelligence tasks. Comprehensive experimental results demonstrate the effectiveness of our proposed approach.

## 2 MOTIVATION

For code intelligence tasks, existing ICL methods [3, 16] use BM25 [45] as the retrieval engine to retrieve examples. The upper left part of Figure 3 presents two test samples collected from the Conala dataset of the program synthesis task. The upper right part shows a

training sample for each test sample. For the first case, the queries of the test sample and the training sample, i.e., Query 1 and Input 1, differ significantly in the used tokens. Therefore, BM25 will not retrieve this training sample for Query 1. However, the structure and functionality of their outputs, i.e., Ground Truth 1 and Output 1, are very similar, both sending a signal to a specified process. So this training sample can be a valuable demonstration for LLMs. For the second case, Query 2 and Input 2 are lexically similar, so BM25 tends to retrieve this training sample as a demonstration. However, they are semantically different, resulting in Ground Truth 2 and Output 2 differ significantly. Using such a demonstration for ICL may mislead LLMs. These two cases show that domain knowledge and semantic understanding ability are important when retrieving demonstrations.

| Samples from the Conala test set | Samples from the Conala training set |
| --- | --- |
| **Query 1:** send a signal `signal.SIGUSR1` to the current process | **Input 1:** kill a process with id `process.pid` |
| **Ground Truth 1:**<br>os . kill ( os . getpid ( ) , signal . SIGUSR1 ) | **Output 1:**<br>os . kill ( process . pid , signal . SIGKILL ) |
| **Query 2:** decode a hex string '4a4b4c' to UTF-8 | **Input 2:** add unicode string '1' to UTF-8 decoded string '\xc2\xa3' |
| **Ground Truth 2:**<br>bytes . fromhex ( '4a4b4c' ) . decode ( 'utf-8' ) | **Output 2:**<br>print ( '\xc2\xa3' . decode ( 'utf8' ) + '1' ) |
| **A sample from the CSN-Java dataset** | |

| Query 3: |
| --- |
| ...<br>if (source instanceof Observable) {<br>    return RxJavaPlugins.onAssembly((Observable<T>) source);<br>}<br>return RxJavaPlugins.onAssembly(new ObservableFromUnsafeSource<T>(source));<br>... |
| **Ground Truth 3:**<br>Wraps an ObservableSource into an Observable if not already an Observable |
| **Input 3:**<br>...<br>if (source instanceof Callable) {<br>    return new RxInstrumentedCallableObservable<>(source, instrumentations);<br>}<br>return new RxInstrumentedObservable<>(source, instrumentations);<br>... |
| **Output 3:**<br>Wrap a observable |

**Figure 3: Examples that BM25 will fail in such cases.**

The lower part of Figure 3 shows a test sample and a training sample collected from the CSN-Java dataset of the code summarization task. Due to space limitation, we omit some code snippets in Query 3 and Input 3. We can see that the input code snippets of the two samples, i.e., Query 3 and Input 3, share similar code structures and their code summaries, i.e., Ground Truth 3 and Output 3, are also alike. This training sample may help LLM generate a good summary for Query 3. However, BM25 overlooks this example since it does not take structural information into account.

To address these issues, we propose the Instructive Code Retriever (ICR). We use contrastive learning to train the ICR based on the feedback from an LLM, leveraging the semantic understanding ability of the LLM to better analyze the input. To incorporate the structural information, we introduce a tree-based loss computation approach which is applicable to code tasks. Details of ICR will be described in the following section.

# 3 INSTRUCTIVE CODE RETRIEVER

## 3.1 Overview

In this section, we will outline the training and inference process of ICR. The retriever aims to retrieve examples from the training set that help the large language model in the inference process.

As illustrated in Figure 2, the origin input is some sample from the training set. In preprocessing and preparation, we first concat task-specific instructions to the origin input. Subsequently, we construct syntax trees for both program and natural language. After preprocessing, we get the input form shown in the lower left part of Figure 2 as a query. For each query, the first step is to retrieve a set of relevant examples from the corpus. After that, we use a large language model to score whether examples would be helpful for LLM inference. Based on the feedback from the language model, we use contrastive learning to learn patterns of the examples that are more useful for the given task. We continuously optimize the training data and iteratively train the retriever. After training, ICR can identify which example is helpful for LLM inference and which is not for a given sample. Once the retriever is trained, during inference, for each test set sample, we retrieve examples most beneficial for the large language model inference process. Finally, we use this sample and retrieved examples as input to the LLM for inference.

## 3.2 Preprocessing and Preparation

In this section, we will describe the preparatory process required for training ICR.

**Building Instructions of Tasks.** ICR is a unified framework for multiple tasks. To enable the model to distinguish between different tasks, we need to incorporate task-specific instructions into the query. Task instructions are short texts describing the task objectives. For example, in the code summarization task, a task instruction could be "Comment on the code." We concatenate the task instruction with the original input for model training, forming the ICR model query in the following format: $(x, y) = T_i \oplus (m, n)$, where $T_i$ is the task instruction, $(m, n)$ is the input-output data pair from training datasets, and $\oplus$ is the concatenation operator.

**Building Syntax Trees.** Traditional bi-encoder architectures typically employ encoders $E_q$ and $E_d$ to encode input and example separately [24], followed by the computation of their similarity scores. The objective of training ICR is to distinguish which examples are highly relevant to the query. The most intuitive approach is to encode both the query and the examples directly and then compute the similarity. However, programming languages inherently possess rich structural information, which is as important as textual information for distinguishing code [60]. Considering that natural language $T_i$ can also be represented as a syntactic tree, the contribution of structural information to the retriever's effectiveness may be augmented [4]. Therefore, we introduce the syntactic tree structure similarity to the traditional bi-encoder architecture to assist the retriever in better understanding the similarity between queries and examples.

The query and examples include both programming languages and natural languages e.g., the input $m$ of program synthesis tasks is a natural language description, and instructions are also in natural language. In our approach, we adopt the 'javalang' [13] and 'ast' [15] libraries to derive the abstract syntax tree of code, while for the
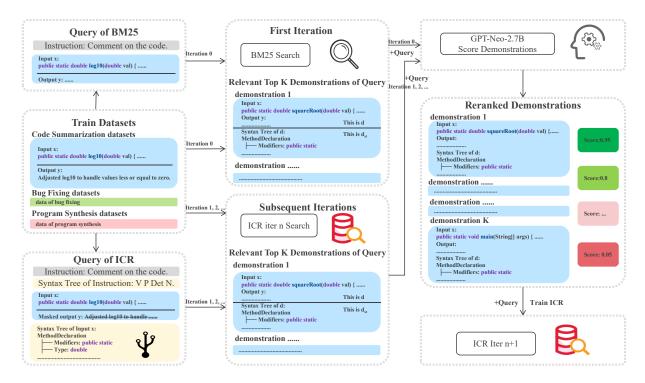
**Figure 2: The Training Workflow of ICR.**

natural language query, we rely on the 'spacy' [20] library to extract its syntax tree. Finally, we get the preorder traversal of the syntax tree of query $x_{st}$ and the syntax tree of example $d_{st}$. The extracted syntax trees will be used in several places, e.g. the model training, the similarity calculation between the query and the example, and the calculation of the loss function. We will describe the details of these processes in the following sections.

**Selecting Relevant Examples for Each Sample.** The primary objective of ICL is to procure high-quality examples for reference by LLM. To achieve this, the retriever plays a pivotal role in selecting exemplary instances. Since the relevant examples found in the previous section may not aid in model inference, we need to use contrastive learning to enable the retriever to identify helpful examples. We need to score examples of each sample for contrastive learning. However, the computational expense associated with scoring the entire training set scales quadratically with its size, rendering such a task both impractical and redundant. Hence, our preliminary step is to roughly screen to find some examples that are similar to the given sample. For each query from the training set, we use only the top K examples (except the query itself) initially retrieved as training data, enabling the model to learn patterns of examples that are useful for LLM inference. This reduces unnecessary computational burden while ensuring the acquisition of high-quality instances for the training process.

Specifically, we divide our process of selecting relevant examples for each sample into the following two scenarios:

- **BM25:** In iter 0, when our retriever is not yet trained, we search for the top K examples (i.e., input-output pair $d$) with the highest BM25 scores for each sample in the training set.

Noting that we are currently processing the training set and can fully utilize ground truth knowledge, to narrow down the scope and select the most relevant examples for a sample, we calculate the similarity between the query of iter 0 (i.e., $(x, y)$) and other examples in the training set.

- **ICR:** In iter 1, 2, ..., n, once our retriever is trained, we iteratively use the trained retriever to find the top K examples (i.e., input-output pair $d$) with the highest scores for each masked query $(x, -) = T_i \oplus (m, -)$ in the training set. At this stage, the ground truth $n$ of the query is masked to ensure consistency between training and testing phases. The inputs for ICR include the query itself and its syntax tree, as well as the examples and their syntax trees. Figure 4 shows the model architecture, and the training process will be introduced in 3.3.

**Scoring for Relevant Examples.** We aim to train ICR to distinguish how useful each relevant example is for handling the given query. To achieve this goal, we need a reference ranking of the relevant examples mentioned above to guide the training. We propose to leverage LLMs to obtain such reference ranking. Specifically, for a given sample $(x, y)$ and its set of examples $(d_1, d_2, ..., d_k)$, we define the score of an example as the probability that the LLM $G$ outputs the ground truth result $y$ based on the input $x$ and the example $d_i$, as follows:

$$S(d_i) = P_G(y \mid d_i, x) \tag{1}$$

Based on the scores of examples, we obtain the ranking $r(d_i)$ of examples, which will guide the training of the retriever in subsequent sections. The function for $r(d_i)$ is as follows:

$$r(d_i) = \text{rank}\left(j < i \mid S(d_j) > S(d_i)\right) \quad (2)$$

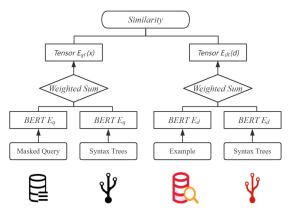Examples with higher scores are ranked higher and have smaller $r(d_i)$.



**Figure 4: The Architecture of ICR.**

## 3.3 Training ICR

ICR is responsible to calculate the similarity or relevance between a masked query and an example. Figure 4 illustrates the architecture of ICR. Given a masked query $x$ an example $d$, ICR first encodes $x$ and its syntax tree $x_{st}$ using a BERT encoder $E_q$, and encodes $d$ and its syntax tree $d_{st}$ using another BERT encoder $E_d$. Note that $d$ contains both the input and the output of the example, and each BERT encoder outputs the contextual embedding of the 'CLS' token as the embedding of the inputs. Then, we aggregate the embeddings of $x$ and $x_{st}$ and the embeddings of $d$ and $d_{st}$, respectively, as follows:

$$E_{qt}(x) = \alpha_1 E_q(x) + \beta_1 E_q(x_{st}) \quad (3)$$
$$E_{dt}(d) = \alpha_2 E_d(d) + \beta_2 E_d(d_{st}) \quad (4)$$

where $\alpha_1$, $\beta_1$, $\alpha_2$, $\beta_2$ are four trainable parameters. Finally, we regard the dot product between the vectors $E_{qt}(x)$ and $E_{dt}(d)$ as the similarity score between the query and the example, as follows:

$$\text{sim}_{\text{tree}}(x, d) = E_{qt}(x)^\top E_{dt}(d) \quad (5)$$

The loss function for training consists of two parts:

**In-Batch-Tree Loss:** To enable the retriever to distinguish the best examples from others, we get inspiration from the loss function of Karpukhin et al. [24]. We modify the in-batch loss function by adding structural information to ensure it aligns with the requirements of code-related tasks. We take into account structural similarity by using the similarity function defined above. We name the loss function as the in-batch-tree loss. The definition of the in-batch-tree loss function is as follows:

$$\mathcal{L}_{bt} = -\log \frac{e^{\text{sim}_{\text{tree}}(x, d^+)}}{e^{\text{sim}_{\text{tree}}(x, d^+)} + \sum_{d^- \in \mathbb{Z}} e^{\text{sim}_{\text{tree}}(x, d^-)}} \quad (6)$$

Where the positive example $d^+$ denotes the highest-scoring example of $x$ within a batch, while $\mathbb{Z}$ encompasses all other examples of $x$ and examples of other queries within this batch. Notice that the negative sample contains not only suboptimal examples of the current query

but also examples of other queries in the same batch (irrelevant examples). The in-batch-tree loss function adopts the principle of contrastive learning, empowering the retriever to pinpoint the example within a batch that offers maximal assistance in inferring the ground truth.

**Ranking-Tree Loss:** The in-batch-tree loss facilitates the retriever's discernment between the highest-quality examples and their suboptimal counterparts. However, the ordering among suboptimal examples also warrants consideration. Therefore, drawing inspiration from Burges et al. [9] and Li et al. [27], we introduce a ranking-tree loss function, defined as follows:

$$\mathcal{L}_{rt} = \sum_{d_i, d_j \in D} r * \log\left(1 + e^{\text{sim}_{\text{tree}}(x, d_j) - \text{sim}_{\text{tree}}(x, d_i)}\right) \quad (7)$$

$$r = \max\left(0, \frac{1}{r(d_i)} - \frac{1}{r(d_j)}\right) \quad (8)$$

This loss function optimizes $sim(x, d_i) > sim(x, d_j)$ when the ranking of $d_i$ is higher than that of $d_j$, i.e., $r(d_i) < r(d_j)$. Consequently, the retriever calculates higher similarity scores for examples ranked higher. The greater the disparity in rankings $r$, the larger the optimization, thus enlarging the gap in similarity between $d_i$ and $d_j$.

The overall tree-based loss function is as follows:

$$\mathcal{L} = \gamma_1 * \mathcal{L}_{bt} + \gamma_2 * \mathcal{L}_{rt} \quad (9)$$

where $\gamma_1$ and $\gamma_2$ are hyper-parameters.

In summary, our training objective is to enable the retriever to identify which examples provide the most assistance in obtaining ground truth from the input in the training set. We achieve this by obtaining the ranking of example quality through feedback from the large model on the training set examples. We train the retriever based on this ranking.

## 3.4 Inferencing for Specific Tasks

During inference, we pre-encode all training set examples using $E_{dt}$. This is because our retriever relies on dot product similarity to rank examples. Precomputing $E_{dt}$ helps avoid redundant computations. For each test example, we get its masked query $(x_{test}, -)$ and compute its encoding $E_{qt}$, and then use the FAISS [23] library to search for a series of examples in the training data that maximize the inner product $sim_{tree}(x_{test}, d)$. We then sort the examples from highest to lowest, resulting in a sequence $D = (d_1, d_2...d_L)$. The final input prompt $P$ provided to the model depends on the maximum input length $|C_i|$ and the maximum output length $|C_o|$ of the model, ensuring that $|P| = \sum_{i=1}^{L} |d_i| + ... + |x_{test}| + |C_o| < |C_m|$. We set a hyperparameter $|C_m|$ and ensure $|C_m| < |C_i|$, instead of using a fixed number of examples. This effectively avoids the example truncation issue. and such that $P$ is the maximal length among all $P_i$ satisfying the condition. We place examples with higher similarity closer to the input. The input prompt format for LLM is as follows:

$$Instruction \parallel Exemplar_1 \parallel \cdots \parallel Input \Rightarrow Prompt \quad (10)$$

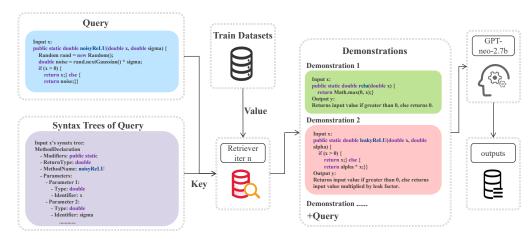With this prompt, we can input it into the large model and obtain the final output.

**Figure 5: The Inference Workflow of ICR.**

## 4 EXPERIMENTAL DESIGN

In this section, we describe the experimental setup that we follow to evaluate the effectiveness of our ICR. Our experiments aim to answer the following research questions:

**RQ1: How effective is ICR enhance the LLM compared to existing state-of-the-art methods?** The primary goal of In-context Learning is to enhance the performance of language models in solving specific tasks. Therefore, we selected two language models of different sizes and conducted comparative experiments to investigate the effectiveness of ICR. In this section, we compared ICR with state-of-the-art In-context Learning approaches for code intelligence tasks: BM25 [16] and *ASAP* [3].

**RQ2: What are the impacts of different modules in the training of ICR?** In RQ2, we will individually dissect the impact of two components, i.e. structural information and iterative training based on LLM feedback, on the performance of our model, serving as ablation studies.

**RQ3: Can our ICR achieve satisfactory performance on un-trained code datasets compared with BM25 and ASAP [3]?** In RQ3, we will compare the performance of JavaScript, PHP and Ruby datasets, which are not used during ICR training, to demonstrate the generalization ability of our model on untrained datasets.

**Table 1: Statistics of the benchmark datasets.**

| | Datasets | Train | Validation | Test |
|---|---|---|---|---|
| Code Summarization | CSN-Java | 164923 | 5183 | 10955 |
| | CSN-Python | 251820 | 13914 | 14918 |
| | CSN-JavaScript | 58025 | 3885 | 3291 |
| | CSN-PHP | 241241 | 12982 | 14014 |
| | CSN-Ruby | 24927 | 1400 | 1261 |
| Bug Fixing | $B2F_{small}$ | 46680 | 5835 | 5835 |
| | $B2F_{medium}$ | 52364 | 6546 | 6545 |
| Program Synthesis | Conala | 2379 | - | 500 |

### 4.1 Datasets

Our experiment involves three code intelligence tasks. Table 1 presents the statistics of the datasets we used in our experiments.

**Bug Fixing:** The bug-fixing dataset used in our study is the popular B2F dataset collected by Tufano et al. [51] and included in

the CODEXGLUE benchmark [35]. Following the original paper's setup, the B2F dataset is divided based on the length of code tokens into two parts: $B2F_{medium}$ and $B2F_{small}$. This dataset was also used by Gao et al. to evaluate the enhancement of LLM through ICL. Hence, we employ the B2F dataset for comparison with our proposed approach. The dataset consists of pairs of buggy code snippets and their corresponding fixed versions in Java, where each pair represents a bug-fixing commit extracted from GitHub repositories. We directly adopt the original dataset splits for our experiments.

**Program Synthesis:** For the program synthesis task, we utilize the CoNaLa dataset [56]. The CoNaLa dataset, short for Code-/Natural Language (CoNaLa), is a widely used benchmark dataset designed for the task of program synthesis from natural language descriptions [16, 37, 61]. This dataset consists of Python data collected from Stack Overflow. With its diverse range of programming scenarios, varying levels of complexity, and high-quality annotations, the CoNaLa dataset serves as a worth considering dataset. We directly adopt the original dataset splits for our experiments.

**Code Summarization:** For the code summarization task, we employ the CodeSearchNet dataset [22]. This dataset contains millions of functions from open-source projects across multiple programming languages, making it popular for training and evaluating models on code summarization tasks. Due to its high quality, the CodeSearchNet dataset has also been included in CODEXGLUE. Code summarization involves generating concise and accurate summaries of code functionality, comments, or documentation. This dataset provides a diverse range of code examples, enabling the development of models that can effectively summarize code across different programming languages and domains. We directly adopt the original dataset splits for our experiments.

### 4.2 Metrics

Since we evaluate the performance of ICR on multiple tasks that have different outputs, we used different metrics as follows.

**BLEU:** BLEU [39] stands as a widely used metric in assessing bug fixing and code summarization tasks. It quantifies the resemblance between automatically generated text and reference text. BLEU evaluates this resemblance by analyzing the overlap of n-grams

between the candidate text and multiple reference texts. Precision scores for n-grams in the generated text are computed, weighted by their frequency.

**ROUGE-L:** ROUGE-L [30] is a metric used to evaluate the quality of summaries by measuring the longest common subsequence (LCS) between the generated summary and a reference summary. ROUGE-L captures the sentence-level structure similarity by considering the sequence of words. This makes it particularly effective for assessing how well a summary retains the order and context of the information from the source text, providing a more nuanced understanding of its coherence and relevance.

**chrF** chrF [40] is a character-level evaluation metric, which measures the quality of outputs by comparing character n-grams between the system output and reference sentences. By operating at the character level, chrF effectively captures subtle differences in morphology, spelling, and word inflections, making it particularly valuable for evaluating translations in morphologically rich languages or low-resource scenarios. This metric calculates an F-score based on the precision and recall of character n-gram overlaps.

**METEOR** METEOR [46] is an evaluation metric designed to improve traditional metrics like BLEU by incorporating linguistic features such as stemming and synonymy. It aligns words between the system output and reference sentences, allowing for partial matches and better handling of varied expressions of meaning. METEOR balances precision and recall, with a slight preference for recall, to reward outputs that capture more reference content. Additionally, it imposes penalties for differences in word order, making it more sensitive to fluency and grammaticality.

**CodeBLEU:** CodeBLEU [43], an adaptation of the BLEU metric, plays a crucial role in evaluating code generation models. By integrating abstract syntax trees, CodeBLEU enhances its evaluation capabilities, infusing code syntax into the assessment process. Furthermore, it leverages data flow to incorporate semantic information, thereby enriching the evaluation process. This integration enables a more comprehensive assessment of code generation models, considering both syntactic correctness and semantic fidelity in the generated code.

**CrystalBLEU** CrystalBLEU [12] is an enhanced version of the traditional BLEU metric designed to improve machine translation evaluation by incorporating linguistic insights. While BLEU focuses on n-gram precision and exact word matches, CrystalBLEU extends this by considering syntactic structures and semantic equivalence, allowing it to account for variations like paraphrases and synonymy. This makes CrystalBLEU more effective in capturing the nuances of translation quality, especially in scenarios where different but equally valid translations may use varied word choices or syntactic structures.

### 4.3 Base Models

In this paper, we employ the ICL method to retrieve examples that enhance the performance of LLMs. In order to show our ICR can work with different LLMs, we selected two popular open-source LLMs. Since the Codex API (Code-Davinci-002), which was widely used in previous research [3, 16], has been deprecated by OpenAI and is no longer available, and other API models may undergo performance changes with version updates or be deprecated, we

prefer to use open-source models to ensure the reproducibility of our work. We respectively enhance each LLM with ICR and other baseline methods. The introductions to these two LLMs are as follows:

**GPT-Neo-2.7B** GPT-Neo-2.7B [7] is an open-source language model developed by EleutherAI as a distilled version of OpenAI's GPT-3. With 2.7 billion parameters, GPT-Neo-2.7B offers substantial performance comparable to much larger models while being more computationally efficient. It is designed to perform a wide range of tasks, demonstrating strong generalization abilities across various domains without fine-tuning. This makes it particularly suitable for tasks requiring significant computational resources and those that benefit from a distilled model's efficiency. We utilize GPT-Neo-2.7B as our foundational model for scoring and inference because it is a distilled model with moderate size. It has good generalization ability and a relatively fast inference speed. Numerous researchers employing GPT-Neo-2.7B as the base model have consistently demonstrated its effectiveness. [27, 48]

**Code Llama-13B** Code Llama-13B [47] is an open-source large language model developed by MetaAI, specifically designed for code generation tasks. It has been widely used in previous research. [5, 32]. With 13 billion parameters, it leverages advanced natural language processing capabilities to generate code based on given text prompts. Code Llama-13B is part of the Llama series of models and aims to provide developers with a powerful tool for various programming tasks, such as code completion, synthesis, and summarization. Its open-source nature allows for extensive customization and adaptation to different coding environments and requirements.

### 4.4 Baselines

In this paper, we will compare our approach with three baseline methods.

**Random** Random represents the random selection of examples, which is typically used as a vanilla baseline model in ICL work [16, 38] to demonstrate the effectiveness of proposed approaches.

**BM25** Best Matching 25 (BM25) [45] is a ranking function used in search engines to evaluate the relevance of documents to a search query. It enhances the classic TF-IDF model by incorporating term frequency saturation and document length normalization. BM25 scores are calculated using a formula that balances term frequency with document length, ensuring that frequently occurring terms and varying document lengths are appropriately weighted. It is widely employed in code-related tasks, due to its effectiveness, simplicity, and robustness. [3, 16, 38] Gao et al. shows that BM25 outperforms pre-trained dense retrievers like UniXcoder [18] and CoCoSoDa [49], and achieves state-of-the-art performance in code intelligence tasks.

**ASAP** Automatic Semantic Prompt Augmentation (ASAP) [3] enhances prompts for code summarization by incorporating three types of designed semantic facts: repository name and path, tagged identifiers, and data flow graphs (DFGs). The ASAP pipeline involves configuring a LLM, a pool of exemplars, and a static analysis tool. The process starts by using BM25 to retrieve relevant exemplars, which are then analyzed to extract designed semantic facts. These facts are included in the prompts to LLM, which generates the final output. This method significantly improves performance

on the code summarization task by providing LLM with detailed and relevant context.

## 4.5 Implementation Details

We trained ICR on $B2F_{small}$, $B2F_{medium}$, CSN-Java, CSN-Python, and Conala datasets. As with all ICL methodologies, we abstain from updating the large model, focusing solely on retriever training. Due to computational resource constraints, we trained ICR using only a quarter of the CSN-Python dataset, and we trained CSN-Java in four slices to accelerate the training process. In the process of selecting relevant examples for each sample, we opt for K=50 relevant examples per sample. Due to computational limitations, we conducted three rounds of iterative training, and we observed that the loss stopped decreasing after four epochs, so we trained ICR for four epochs in each iteration. Training takes place on a server equipped with 8 NVIDIA A800-80G GPUs. To address the repetition problem, we implemented post-processing to trim the repetitive parts, retaining only distinct instances. Moreover, because we need to train on multiple datasets with significant size differences, the sampling rate for Conala is set to 3 per epoch, while for the other datasets, it is set to 0.7. $|C_m|$ is set to 2048, which is the maximum input length for GPT-Neo-2.7B. Considering that the in-batch-tree loss calculates examples within a batch of queries and ranking-tree loss only calculates examples of the query itself, the in-batch-tree loss is generally larger than the ranking-tree loss. We set $\gamma_1$ to 1 and $\gamma_2$ to 4 to ensure both loss functions are effective. In our implementation, we fixed the random seed for consistent results, while employing a greedy decoding strategy when invoking the LLM.

## 5 RESULT ANALYSIS

### 5.1 Effectiveness of ICR

In RQ1, we compare ICR with BM25 on the three tasks. Because ASAP only reports its performance on code summarization, we only compare ICR with ASAP in the two code summarization datasets. For our base model selection, we first compare the performance of ICR with baseline models on the open-source GPT-Neo-2.7B [7]. Table 2 and 3 presents the evaluation results. We can observe that on the code summarization task, ICR model significantly outperforms BM25, achieving BLEU-4 score increases of 136.2% on the CSE-Java dataset and 137.5% on the CSN-Python dataset. Additionally, it improves over BM25 in ROUGE-L by of 70.6% and 68.2% on these datasets, respectively. Compared to ASAP, a method specifically designed for code summarization, ICR demonstrates improvements of 50.0% on the CSN-Java dataset and 90% on the CSN-Python dataset in terms of BLEU-4. Additionally, ICR achieves increases in ROUGE-L scores of 34.3% and 60.8% on these respective datasets. ICR also improves chrF score of over 30% and METEOR over 41.3%.

For the program synthesis task, ICR outperforms the BM25 method by 74.6% on CodeBLEU and 67.2% on CrystalBLEU. Conala is a dataset with a relatively small training set, and when the training set is insufficient, ICR can learn about the quality of examples from other datasets of different tasks. In contrast, the small dataset limits the comprehensiveness of the corpus that BM25 can retrieve. In the Bug Fixing task, ICR outperforms the BM25 method by 3.6

BLEU-4 score on $B2F_{small}$ and 3.2 BLEU-4 score on $B2F_{medium}$, respectively.

To evaluate the generality of our retriever, we transferred ICR trained with feedback from GPT-Neo-2.7B to Code Llama-13B. As shown in Table 2 and 3, our model significantly outperforms the baseline methods. For the code summarization task, ICR outperforms the BM25 method by 32.4% on the CSN-Java dataset and 34.1% on the CSN-Python dataset in terms of BLEU-4, while improving ROUGE-L by 19.2% and 15.6%. Further, it also surpasses the heuristic-based ASAP method by 16.8% and 18.3% on BLEU-4, while improving ROUGE-L by 13.8% and 19.1%. Additionally, it achieved the highest scores on chrF and METEOR. In the program synthesis task, we exceed the BM25 method by 19.7%. In the bug fixing task, ICR's performance also surpasses that of BM25 in BLEU-4 scores. This demonstrates the effectiveness of our ICR when transferred to a more powerful LLM, such as Code Llama-13B.

> When tested with GPT-Neo-2.7B and Code Llama-13B, our method outperforms BM25 and ASAP in all three code intelligence tasks.

### 5.2 Ablation Study

To evaluate the effectiveness of each component of ICR, we conducted ablation studies on the structural information and iterative training based on LLM feedback. The experimental results are shown in the Table 4. We have removed two key modules from ICR and constructed the following two variants:

- **ICR + GPT-Neo-2.7B - w/o Structural Info.**: The ablation of structural information involves whether the structural information is used in calculating similarity and loss. Specifically, we remove the tree structure information from our model and replace it all with zero vectors.
- **ICR + GPT-Neo-2.7B - w/o Feedback-driven Training**: The feedback-driven training ablation experiment compares single-iteration training and iterative training based on LLM feedback. When building variant of single-iteration, ICR only using BM25 to initialize the examples in the training phase. Therefore, this variant is unable to iteratively train based on the example scores provided by LLM feedback.

*5.2.1 Structural Information.* Compared to the variant that does not incorporate structural information, ICR demonstrated an improvement of 22.0% on CSN-Java and 15.2% on CSN-Python in terms of BLEU-4 for the code summarization task. Additionally, CodeBLEU improves 38.6% on the Conala dataset for the program synthesis task. In the bug fixing task, structural information improves 1.3 BLEU-4 points on $B2F_{small}$. Syntax trees contain rich structural information [6], which can provide valuable cues for code intelligence tasks. In addition to characterizing structural information during model input, our method also calculates structural information during training. The loss function we designed includes structural information that can effectively help the model understand the input structure and retrieve structurally more similar samples.

*5.2.2 Feedback-Driven Iterative Training.* Feedback-driven iterative training has led to significant performance improvements,

**Table 2: Performance of ICR on Code Summarization.**

| Tasks | Code Summarization | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Datasets | CSN-Java | | | | CSN-Python | | | |
| Metrics | ROUGE-L | BLEU-4 | chrF | METEOR | ROUGE-L | BLEU-4 | chrF | METEOR |
| Random + GPT-Neo-2.7B | 18.5 | 3.6 | 16.8 | 17.1 | 18.8 | 2.8 | 15.6 | 16.0 |
| BM25 + GPT-Neo-2.7B | 20.4 | 4.7 | 19.0 | 18.3 | 19.5 | 3.2 | 16.5 | 16.3 |
| ASAP + GPT-Neo-2.7B | 25.9 | 7.4 | 23.0 | 22.5 | 20.4 | 4.0 | 17.7 | 16.6 |
| ICR + GPT-Neo-2.7B | **34.8** | **11.1** | **29.9** | **31.8** | **32.8** | **7.6** | **26.4** | **28.4** |
| Random + Code Llama-13B | 34.0 | 8.1 | 27.7 | 28.3 | 32.7 | 6.4 | 25.8 | 27.1 |
| BM25 + Code Llama-13B | 33.8 | 10.5 | 29.4 | 29.5 | 34.0 | 8.2 | 28.4 | 28.3 |
| ASAP + Code Llama-13B | 35.4 | 11.9 | 29.1 | 30.7 | 33.0 | 9.3 | 27.1 | 27.9 |
| ICR + Code Llama-13B | **40.3** | **13.8** | **33.0** | **36.1** | **39.3** | **11.0** | **30.9** | **34.1** |

**Table 3: Performance of ICR on Program Synthesis and bug fixing.**

| Tasks | Program Synthesis | | | Bug Fixing | | | |
|---|---|---|---|---|---|---|---|
| Datasets | Conala | | | B2F_small | | B2F_medium | |
| Metrics | CodeBLEU | BLEU-4 | CrystalBLEU | BLEU-4 | CrystalBLEU | BLEU-4 | CrystalBLEU |
| Random + GPT-Neo-2.7B | 16.6 | 15.4 | 11.5 | 75.0 | 41.7 | 87.1 | 48.7 |
| BM25 + GPT-Neo-2.7B | 18.5 | 17.6 | 12.8 | 75.6 | 41.9 | 86.8 | 48.7 |
| ICR + GPT-Neo-2.7B | **32.3** | **31.6** | **21.4** | **79.2** | **42.1** | **89.0** | **49.1** |
| Random + Code Llama-13B | 36.1 | 35.4 | 25.1 | 78.2 | 42.0 | 88.7 | 48.7 |
| BM25 + Code Llama-13B | 37.1 | 38.0 | 26.4 | 78.8 | 42.6 | 89.8 | 49.3 |
| ICR + Code Llama-13B | **44.4** | **45.8** | **30.4** | **80.1** | **43.5** | **90.7** | **49.7** |

**Table 4: Ablation Studies of ICR.**

| Tasks | Code Summarization | | | | Program Synthesis | | Bug Fixing | |
|---|---|---|---|---|---|---|---|---|
| Datasets | CSN-Java | | CSN-Python | | Conala | | B2F_small | B2F_medium |
| Metrics | ROUGE-L | BLEU-4 | ROUGE-L | BLEU-4 | CodeBELU | BLEU-4 | BLEU-4 | BLEU-4 |
| ICR + GPT-Neo-2.7B | **34.8** | **11.1** | **32.8** | **7.6** | **32.3** | **31.6** | **79.2** | **89.0** |
| - w/o Structural Info. | 29.9 | 9.1 | 30.0 | 6.6 | 23.3 | 22.6 | 77.9 | **89.0** |
| - w/o Feedback-driven Training | 30.2 | 7.7 | 21.7 | 3.0 | 32.1 | 28.7 | 65.8 | 88.6 |

particularly on the CSN-Java (26.0% BLEU-4), CSN-Python (153.3% BLEU-4), and $B2F_{small}$ (20.4% BLEU-4) datasets. For **ICR + GPT-Neo-2.7B - w/o Feedback-driven Training** variant, given a training sample, it relies solely on the retrieval results from BM25 to select relevant examples for training. This approach could be insufficient and suboptimal. For ICR + GPT-Neo-2.7B method, we utilize the ICR model trained in the previous iteration for further retrieval. Some of retrieved examples might still not be helpful for model inference. This distinction is made through feedback from the large model, training the retriever to discard these unhelpful examples in the next round. We refine the model based on the feedback from LLM to enhance its performance. This process represents a effectively continual iteration of corrections and optimizations to ICR. As the iterations progressed, the performance of ICR improved, enabling the acquisition of higher-quality samples and leading to better overall performance.

> For all three code intelligence tasks, both structural information and feedback-driven training helps ICR find higher quality examples.

## 5.3 Transferability of ICR on Untrained Datasets

In RQ3, we explore whether ICR can generalize to languages it has not been trained on. Considering the CSN dataset includes data from other programming languages, we conducted the experiments on the code summarization task. We will compare ICR-enhanced GPT-Neo-2.7B with two baselines: BM25 and ASAP.

To answer RQ3, we evaluate the performance of the ICR approach in comparison with baseline models for retrieving examples from untrained programming languages, specifically on the CSN-JavaScript, CSN-PHP, and CSN-Ruby datasets. We employ ICR that was trained on the CSN-Java, CSN-Python, Conala, $B2F_{small}$, and $B2F_{medium}$ datasets to retrieve similar samples from these untrained datasets. For inference, we use GPT-Neo-2.7B, maintaining all settings consistent with those outlined in Section 3.4 and Section 4.5.

While ICR was only trained on Java and Python datasets during training, the results in Table 5 indicate a performance improvement of 48.5%, 22.5%, and 122.7% with ICR compared to BM25 on CSN-JavaScript, CSN-PHP, and CSN-Ruby in terms of BLEU-4, respectively. Compared to ASAP, our approach showed improvements

**Table 5: Performance of ICR on Untrained Datasets.**

|  | CSN-JavaScript | | CSN-PHP | | CSN-Ruby | |
| --- | --- | --- | --- | --- | --- | --- |
| Metrics | ROUGE-L | BLEU-4 | ROUGE-L | BLEU-4 | ROUGE-L | BLEU-4 |
| BM25 + GPT-Neo-2.7B | 18.5 | 6.8 | 31.9 | 10.2 | 18.3 | 2.2 |
| ASAP + GPT-Neo-2.7B | 22.8 | 9.0 | 36.1 | 10.7 | 25.1 | 3.9 |
| ICR + GPT-Neo-2.7B | 23.8 | 10.1 | 37.6 | 12.5 | 27.7 | 4.9 |

**Table 6: Analysis of Example Order.**

|  | $B2F_{small}$ | $B2F_{medium}$ | Conala | CSN-Java | CSN-Python |
| --- | --- | --- | --- | --- | --- |
| Metrics | BLEU-4 | BLEU-4 | CodeBLEU | BLEU-4 | BLEU-4 |
| ICR Random + GPT-Neo-2.7B | 78.9 | 88.6 | **32.6** | 11.1 | 7.5 |
| ICR Similarity + GPT-Neo-2.7B | **79.2** | **89.0** | 32.3 | 11.1 | **7.6** |
| ICR Reverse Similarity + GPT-Neo-2.7B | 78.8 | 88.3 | 32.4 | **11.2** | 7.3 |

of 12.2%, 16.8%, and 25.6% in terms of BLEU-4, respectively. This suggests that the knowledge acquired by ICR from training data and structural information can be generalized to untrained languages, yielding better performance. This is because ICR effectively learns the structural knowledge of code, which is transferable across languages. Additionally, it extracts deeper insights from the feedback of the large language model. Unlike BM25, which relies solely on text-based example selection, ICR could provide more relevant and insightful examples. This ultimately offers better inspiration for the large language model.

> Our approach also shows superior performance on datasets it has not been trained on.

## 6 DISCUSSION

### 6.1 Findings and Limitations

Gao et al. [16] find that the example selection, quantity, and order are pivotal factors influencing the performance of ICL. While the impact of example selection and example quantity on model performance is widely acknowledged [31, 57], the effects of order need further discussion. Gao et al. observed that while the performance of large models can be affected by the order of examples, this impact diminishes when the quality of examples is high, such as those retrieved using BM25. For instance, in their study focusing on code summarization enhanced by BM25, randomizing the order of examples resulted in optimal performance. In general, high-quality examples show less sensitivity to their order [28, 36], this phenomenon that has been confirmed across a range of natural language processing tasks.

To evaluate the effect of example order, we compare random order and two simple ordering methods, i.e., ICR similarity order and ICR reverse similarity order. The ICR similarity method organizes examples based on their retrieval order, as described in Section 3.4, placing those with higher similarity closer to the input query. Conversely, the ICR reverse similarity order arranges examples so that those with lower similarity are positioned nearer to the input query. As illustrated in Table 6, reversing the order of ICR examples fed into the model has no substantial impact on the results for most datasets. This finding indirectly indicates the high quality of the examples retrieved by ICR and consistently provides valuable

insights into the model, making the order of examples relatively unimportant.

We posit that the fundamental factors impacting the augmentation of large models through ICL lie in example selection and composition. While each individually chosen example may be optimal, their collective efficacy may be compromised due to potential redundancy. Hence, opting for a more varied selection of examples may enhance the inference capabilities of large models. However, an exhaustive examination of all possible example combinations is computationally prohibitive, necessitating the use of approximate algorithms in this regard. Our current study does not delve into example composition, leaving this avenue for future investigation.

### 6.2 Threats to validity

**Threats to internal validity** To enable ICR to perform multiple tasks, we employed different instructions for each task. The variation in instruction selection may impact model performance. To mitigate this threat, we tested the performance of different instructions on a small-scale subset of the corresponding training set and select the best-performing instruction.

**Threats to external validity** ICR adopts a unified architecture for multitask learning, with training conducted solely on three code-related tasks and datasets in two languages, i.e. Python and Java. Performance on the datasets in other languages and in other tasks remains to be tested in future evaluations. To mitigate this threat, we assessed the performance of ICR on datasets in other languages that were not part of its training data (see Section 5.3). Another threat is that Conala is a relatively small data set. This may lead to insufficient training in the program synthesis task. Therefore, it may not be able to retrieve enough high-quality examples. Additionally, this threat could potentially compromise the evaluation of ICR on the code synthesis task, making it less generalizable. In the future, we plan to conduct experiments on more datasets for program synthesis task.

**Threats to construct validity** For the same dataset, we utilized the evaluation metrics identical to those used by Gao et al. [16], and the metric we used is a potential threat, as it may not fully reflect the quality of the output. We did not conduct manual evaluations of ICR, which we plan to undertake in the future.

## 7 RELATED WORKS

### 7.1 Pre-Trained Models for code intelligence tasks

Pre-trained models have achieved significant accomplishments in the field of software engineering. Drawing inspiration from the triumph of large pre-trained models in NLP [8, 11, 33, 50], Feng et al. [14] adapted a pre-trained language model named CodeBERT, leveraging a vast code corpus. Guo et al. [19] further enhanced this approach with GraphCodeBERT, integrating structural code information for superior performance. Both CodeBERT and Graph-CodeBERT have excelled in programming language tasks. Ahmad et al. [2] mirrored BART's [26] model structure, pre-training their PLBART model, tailored for bug fixing. Meanwhile, Wang et al. [53] introduced CodeT5, a pre-trained code model fashioned after T5's [42] architecture, along with specialized pre-training tasks for code. CodeT5 has demonstrated outstanding performance across various downstream tasks, including code generation. These seminal models represent significant strides towards leveraging pre-trained models for code-related tasks, demonstrating remarkable efficacy across diverse downstream applications.

While pre-trained models offer a myriad of benefits, they also come with their fair share of limitations that can hamper their real-world utility. One significant drawback is the mismatch between the pre-training tasks and downstream applications, requiring fine-tuning to bridge the gap. This fine-tuning process can be cumbersome for users, adding complexity to model deployment. Moreover, due to their relatively limited parameter space, pre-trained models often struggle with generalizing across diverse tasks, resulting in a lackluster performance on tasks and datasets that haven't undergone specific fine-tuning procedures.

### 7.2 Large Language Models

With the further development of pre-trained models, models with larger parameters have emerged, known as Large Language Models (LLMs). LLMs demonstrate exceptional generalization capabilities across a wide range of tasks, eliminating the need for fine-tuning to achieve good performance. The GPT series of models developed by OpenAI, including GPT-3 [8], ChatGPT, GPT-4 [1], CodeX [10], and others, have achieved success in various tasks. GPT-Neo [7] is an open-source model distilled by EleutherAI from GPT-3, possessing fewer model parameters yet achieving comparable performance to larger models. Code Llama [47] is an open-source code generation model released by MetaAI, capable of generating code using text prompts.

The advent of large language models has liberated users from the complexities of fine-tuning models. However, these models are particularly sensitive to their input prompts [59]. Hence, the current challenge revolves around creating prompts of superior quality for these large language models. In-Context Learning (ICL) technology steps in to tackle this challenge by furnishing examples to inspire and guide these expansive models.

### 7.3 In-Context Learning

Large Language Models (LLMs) have demonstrated remarkable emergent capabilities, and enhancing them has become a focal point in research. In-context learning (ICL) has shown a series of successes in augmenting large language models. In ICL, the model completes new tasks directly by being provided with a set of contextual examples, without the need for further fine-tuning or parameter modification. These contextual examples typically include input-output pairs, which the model uses to infer the output for new inputs. Brown et al. [8] were among the pioneers to showcase that providing a series of examples can significantly boost the performance of GPT-3. DPR [24], employing a dual-encoder structure and a contrastive learning approach, has achieved notable success in open-domain question-answering tasks. Gao et al. [16] were the first to apply ICL techniques to code-related tasks, utilizing BM25 [45] to retrieve examples and enhance large language models in code intelligence tasks. Ahmed et al. [3] further utilized semantic facts to augment the performance of the code summarization task after BM25 retrieving examples.

However, in the field of code intelligence tasks, existing approaches are still based on lexical similarity retrievers [3, 16]. In this paper, we proposed a novel approach ICR based on contrastive learning to help retrieve highly relevant examples for a given code task. Our approach incorporates both semantic and structural information to help retriever understand the relationship between queries and examples. This retriever will inspire LLMs to generate more accurate results to complete various code intelligence tasks.

## 8 CONCLUSION AND FUTURE WORK

In this paper, we introduced an approach named Instructive Code Retriever (ICR), which is designed to fetch high-quality examples for multiple code intelligence tasks, thereby enhancing the inference of large language models. Given the training data, ICR first uses BM25 (iter 0) or the iteration-trained ICR to score examples and search for relevant ones. After getting the relevant examples, we further leverage LLMs to rank them based on their generation probability. Then we use this ranking to calculate a customized tree-based loss function to achieve the goal of training ICR. During testing, ICR effectively retrieved useful examples to inspire the large model. Across various code intelligence tasks, such as bug fixing, program synthesis, and code summarization, ICR demonstrated significant improvements over baseline models. Additionally, we analyzed each component of ICR and assessed its transferability across datasets in different languages. Overall, our retriever shows promise in assisting researchers and developers across diverse code intelligence tasks and datasets in various languages. In the future, we aim to delve deeper into the impact of example combinations and diversity on inspiring large models, as well as enhancing retrievers' ability to learn code-specific structures. Our replication package is available at https://github.com/kingofheven/ICR.

# REFERENCES

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified pre-training for program understanding and generation. *arXiv preprint arXiv:2103.06333* (2021).

[3] Toufique Ahmed, Kunal Suresh Pai, Premkumar Devanbu, and Earl Barr. 2024. Automatic semantic augmentation of language model prompts (for code summarization). In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.

[4] Shengnan An, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Jian-Guang Lou, and Dongmei Zhang. 2023. How Do In-Context Examples Affect Compositional Generalization? *arXiv preprint arXiv:2305.04835* (2023).

[5] Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q Jiang, Jia Deng, Stella Biderman, and Sean Welleck. 2023. Llemma: An open language model for mathematics. *arXiv preprint arXiv:2310.10631* (2023).

[6] Jiangang Bai, Yujing Wang, Yiren Chen, Yaming Yang, Jing Bai, Jing Yu, and Yunhai Tong. 2021. Syntax-BERT: Improving pre-trained transformers with syntax trees. *arXiv preprint arXiv:2103.04350* (2021).

[7] Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. Gpt-neo: Large scale autoregressive language modeling with mesh-tensorflow. *If you use this software, please cite it using these metadata* 58 (2021), 2.

[8] 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.

[9] Christopher JC Burges. 2010. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11, 23-581 (2010), 81.

[10] 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).

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[12] Aryaz Eghbali and Michael Pradel. 2022. CrystalBLEU: precisely and efficiently measuring the similarity of code. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–12.

[13] João Roberto Ferreira Fabio Zadrozny, Torsten Marek. 2018. javalang: Pure Python Java parser and tools. https://github.com/c2nes/javalang Accessed: 2024-06-07.

[14] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng 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).

[15] Python Software Foundation. 2024. AST: Abstract Syntax Trees. https://docs.python.org/3/library/ast.html Accessed: 2024-06-07.

[16] Shuzheng Gao, Xin-Cheng Wen, Cuiyun Gao, Wenxuan Wang, Hongyu Zhang, and Michael R Lyu. 2023. What makes good in-context demonstrations for code intelligence tasks with llms?. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 761–773.

[17] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. 2017. Program synthesis. *Foundations and Trends® in Programming Languages* 4, 1-2 (2017), 1–119.

[18] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. Unixcoder: Unified cross-modal pre-training for code representation. *arXiv preprint arXiv:2203.03850* (2022).

[19] 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).

[20] Matthew Honnibal and Ines Montani. 2020. spaCy: Industrial-strength Natural Language Processing in Python. https://spacy.io Accessed: 2024-06-07.

[21] Kai Huang, Xiangxin Meng, Jian Zhang, Yang Liu, Wenjie Wang, Shuhao Li, and Yuqing Zhang. 2023. An empirical study on fine-tuning large language models of code for automated program repair. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1162–1174.

[22] 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).

[23] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.

[24] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906* (2020).

[25] Alexander LeClair, Sakib Haque, Lingfei Wu, and Collin McMillan. 2020. Improved code summarization via a graph neural network. In *Proceedings of the 28th international conference on program comprehension*. 184–195.

[26] 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).

[27] Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for in-context learning. *arXiv preprint arXiv:2305.04320* (2023).

[28] Xiaonan Li and Xipeng Qiu. 2023. Finding supporting examples for in-context learning. *arXiv e-prints* (2023), arXiv–2302.

[29] Zhiyu Li, Shuai Lu, Daya Guo, Nan Duan, Shailesh Jannu, Grant Jenks, Deep Majumder, Jared Green, Alexey Svyatkovskiy, Shengyu Fu, et al. 2022. Automating code review activities by large-scale pre-training. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1035–1047.

[30] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.

[31] 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).

[32] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688* (2023).

[33] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

[34] Jiawei Lu, Zhijie Tang, and Zhongxin Liu. 2023. Improving Code Refinement for Code Review Via Input Reconstruction and Ensemble Learning. In *2023 30th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, 161–170.

[35] 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).

[36] 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. *arXiv preprint arXiv:2104.08786* (2021).

[37] Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, et al. 2022. Lila: A unified benchmark for mathematical reasoning. *arXiv preprint arXiv:2210.17517* (2022).

[38] Noor Nashid, Mifta Sintaha, and Ali Mesbah. 2023. Retrieval-based prompt selection for code-related few-shot learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2450–2462.

[39] 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.

[40] Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the tenth workshop on statistical machine translation*. 392–395.

[41] Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2022. Reasoning with language model prompting: A survey. *arXiv preprint arXiv:2212.09597* (2022).

[42] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.

[43] 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).

[44] Peter C Rigby and Christian Bird. 2013. Convergent contemporary software peer review practices. In *Proceedings of the 2013 9th joint meeting on foundations of software engineering*. 202–212.

[45] Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*. Springer, 232–241.

[46] Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. 2021. Reassessing automatic evaluation metrics for code summarization tasks. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1105–1116.

[47] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).

[48] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633* (2021).

[49] Ensheng Shi, Yanlin Wang, Wenchao Gu, Lun Du, Hongyu Zhang, Shi Han, Dongmei Zhang, and Hongbin Sun. 2023. Cocosoda: Effective contrastive learning for code search. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2198–2210.

[50] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[51] Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin White, and Denys Poshyvanyk. 2019. An empirical study on learning bug-fixing patches in the wild via neural machine translation. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 28, 4 (2019), 1–29.

[52] Rosalia Tufano, Luca Pascarella, Michele Tufano, Denys Poshyvanyk, and Gabriele Bavota. 2021. Towards automating code review activities. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 163–174.

[53] 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).

[54] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682* (2022).

[55] Chunqiu Steven Xia, Yifeng Ding, and Lingming Zhang. 2023. The Plastic Surgery Hypothesis in the Era of Large Language Models. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 522–534.

[56] Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. 2018. Learning to mine aligned code and natural language pairs from stack overflow. In *Proceedings of the 15th international conference on mining software repositories*. 476–486.

[57] Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active example selection for in-context learning. *arXiv preprint arXiv:2211.04486* (2022).

[58] Xin Zhout, Kisub Kim, Bowen Xu, Jiakun Liu, DongGyun Han, and David Lo. 2023. The devil is in the tails: How long-tailed code distributions impact large language models. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 40–52.

[59] Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Neil Zhenqiang Gong, Yue Zhang, et al. 2023. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv preprint arXiv:2306.04528* (2023).

[60] Qihao Zhu, Zeyu Sun, Yuan-an Xiao, Wenjie Zhang, Kang Yuan, Yingfei Xiong, and Lu Zhang. 2021. A syntax-guided edit decoder for neural program repair. In *Proceedings of the 29th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*. 341–353.

[61] Terry Yue Zhuo. 2024. ICE-Score: Instructing Large Language Models to Evaluate Code. In *Findings of the Association for Computational Linguistics: EACL 2024*. 2232–2242.