# AdaptiveLLM: A Framework for Selecting Optimal Cost-Efficient LLM for Code-Generation Based on CoT Length

Junhang Cheng, Fang Liu*, Chengru Wu, Li Zhang
State Key Laboratory of Complex & Critical Software Environment, School of Computer Science and Engineering
Beihang University, Beijing, China
chengjunhang7@gmail.com,{fangliu,23230618,lily}@buaa.edu.cn

## Abstract

While Large Language Models (LLMs) have significantly advanced code generation efficiency, they face inherent challenges in balancing performance and inference costs across diverse programming tasks. Dynamically selecting the optimal LLM based on task difficulty and resource constraints offers a promising approach to achieve an optimal balance between efficiency and performance. However, existing model selection methods are resource-intensive and often neglect cost efficiency. Moreover, these approaches rely on human-annotated difficulty labels that are frequently inaccessible in real-world settings and may not align with the LLM's own assessment of task difficulty. In this paper, we introduce AdaptiveLLM, a framework that dynamically selects optimal LLMs for a given coding task by automatically assessing task difficulty. Our framework first estimates task difficulty using Chain-of-Thought lengths generated by reasoning model, clusters these into three difficulty levels via k-means, and fine-tunes CodeBERT to embed difficulty-aware features. A trained XGBoost classifier then selects the best model for each problem, optimizing the performance-cost trade-off. Experimental results show that AdaptiveLLM achieves a 7.86% improvement in pass@1 score while reducing resource consumption by 88.9% compared to baseline method ComplexityNet. When compared to a single model, AdaptiveLLM demonstrates an approximately 15% accuracy improvement, while maintaining the same level of cost consumption. Apart from that, the difficulty assessment using CoT provides more reliable selection criteria than human evaluation. Our replication package is available at https://github.com/cjhCoder7/AdaptiveLLM.

## CCS Concepts

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

## Keywords

Large Language Model, Code Generation, Model Selection
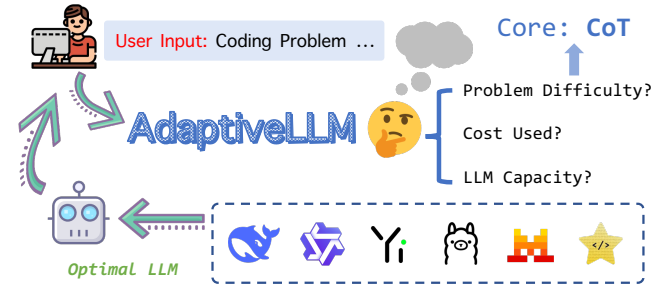
*Corresponding author.

**Figure 1: The Use of AdaptiveLLM.**

## 1 Introduction

Large Language Models have emerged as transformative tools in code understanding and generation, driving significant advancements in programming efficiency through intelligent assistants such as GitHub Copilot [17] and Cursor [4]. These systems leverage LLM's contextual reasoning capabilities to predict and auto-complete code snippets, reducing development time and effort. However, expanding application scenarios reveal two critical challenges: ❶ LLM exhibit varying performance across programming tasks of different complexities, and ❷ their high computational costs for inference and training necessitate optimized resource allocation. Specifically, SWE-Bench [27] has demonstrated that different types of code LLMs possess varying capabilities in solving different types of code problems. Similarly, code LLMs with different parameter sizes also show distinct abilities in addressing code problems of varying difficulties. While larger models generally exhibit stronger problem-solving capabilities, smaller models can achieve comparable results on simpler tasks. Given that larger models incur higher operational costs and are overqualified for simpler problems, *dynamically selecting the optimal LLM based on task complexity and resource constraints presents a promising strategy to balance efficiency and performance.*

To address this challenge, researchers have developed various model selection approaches, which generally fall into two categories. The first category comprises router-based approaches [10, 12, 21, 26, 56]. For example, RouterDC [12] introduces a novel routing mechanism that differs from traditional selection methods by embedding the problem-LLM matching into a vector space. This

approach allows for a single LLM invocation to solve the problem, reducing resource consumption compared to previous strategies that required generating and screening responses from all LLMs in advance. However, RouterDC overlooks the difficulty levels of the problems themselves and relies solely on model responses during the framework construction to create problem-LLM matches, which can be resource-intensive. The second category selects different LLMs based on the difficulty levels of problems [6, 25, 36]. For instance, ComplexityNet [6] uses difficulty tags to guide the selection of optimal models. However, this approach overlooks cost considerations during model invocation, frequently relying on high-performance but expensive closed-source models like GPT-3.5 and GPT-4. Additionally, ComplexityNet relies on annotated difficulty tags, which are often unavailable in real-world settings, and they may not accurately reflect an LLM's intrinsic perception of task difficulty. Furthermore, directly estimating difficulty using LLMs has proven unreliable due to the randomness in their predictions.

Recent advances in reasoning models have opened new possibilities for automating problem difficulty assessment [19, 39, 40, 46]. These models produce step-by-step Chains-of-Thought (CoT) reasoning process, where longer reasoning sequences often correlate with higher problem complexity. This method not only mimics human cognitive patterns but also provides a systematic framework for evaluating difficulty. By integrating reasoning models into the framework and considering model costs, we can potentially overcome the limitations of both Router-based and difficulty-based methods for code generation tasks.

To this end, we propose AdaptiveLLM, a framework that dynamically selects optimal code generation models based on automated difficulty assessment. AdaptiveLLM first estimates the difficulty of each coding problem according to the CoT length generated from LLMs with enhanced reasoning capabilities. These CoT lengths are clustered into three difficulty levels using k-means, and the resulting labels are used to fine-tune the CodeBERT embedding model [16]. This fine-tuning process enriches the problem embeddings with difficulty-aware features. Subsequently, we train an XGBoost classifier to select the best-performing model for each problem, considering both response quality and resource efficiency to generate a balanced ranking. By dynamically matching problems with the most suitable models, AdaptiveLLM optimizes the trade-off between performance and cost. **Experimental results on three datasets of varying difficulty show that AdaptiveLLM achieves a 7.86% improvement in capability and an 88.9% reduction in cost compared to baseline methods.**

Our contributions are summarized as follows:

- We present AdaptiveLLM, a framework for selecting the most cost-effective LLM based on problem difficulty, model capability, and model cost. It enables personalized selection of the optimal LLM for each problem.
- We propose a novel programming task difficulty assessment method based on CoT of reasoning models. This automated approach eliminates the need for human intervention by using the length of CoT to estimate problem complexity.
- We evaluate AdaptiveLLM and the baseline method on benchmark datasets. Results show that AdaptiveLLM achieves superior

performance with significantly lower resource consumption. Additionally, our analysis of the CoT difficulty assessment method reveals that CoT more accurately reflects LLMs' perception of problem difficulty compared to human labels.

## 2 Related Work

### 2.1 LLMs for Code Generation

Large Language Models have achieved significant breakthroughs in natural language processing, particularly demonstrating exceptional capabilities in code understanding and generation [1, 18, 19, 31, 44–46, 48, 52]. Previous works focused on fine-tuning pretrained LLMs to address tasks of varying types and complexities. For example, mathematical reasoning tasks have been addressed through specialized fine-tuning of models such as WizardMath [35], Qwen2.5-Math [53] and MetaMath [55]. Similarly, code generation capabilities have been improved through fine-tuning, as demonstrated by models such as Qwen2.5-Coder [24] and DeepSeek-Coder-V2 [58]. These models exhibit significant performance gains in their respective target domains.

Code generation benchmarks have evolved to assess LLM's capabilities across complexity levels. HumanEval [11] is a popular benchmark, whitch then extended to multilingual adaptations [9, 41, 57], novel code completion paradigms [37], and enhanced testing frameworks [32]. High-complexity evaluation leverages CodeContests [29] for algorithmic challenges and SWE-Bench [27] for real-world software engineering scenarios. FullStackBench [33] further enables cross-domain evaluation spanning 16 programming languages with SandboxFusion execution.

Our study selects three datasets: HumanEval, LeetCodeSample, and CodeContests. Additionally, we choose eight code LLMs and DeepSeek R1 reasoning model.

### 2.2 Evaluation of Coding Problems

The difficulty of programming problems is a crucial factor in selecting the optimal model. Various approaches have been proposed for assessing the difficulty of programming problems.

LeetCodeSample [33] and CodeContests [29] extracted from competitive programming platforms, such as LeetCode [3] and Codeforces [2], often rely on user performance data to quantify problem difficulty. However, these metrics suffer from cross-platform incomparability and dependence on user performance data. Another approach involves analyzing the complexity of solutions. Wang et al. [49] evaluated problem difficulty using five code complexity metrics: Line Complexity, Cyclomatic Complexity [15], Halstead Complexity [22], Cognitive Complexity [8], and Maintainability Index [51]. However, not all the problems have standard solutions. Furthermore, for problems with multiple valid solutions, the variability in code complexity across different implementations will influence it. Bae et al. [6] and Jeong et al. [25] evaluate difficulty based on the number of interactions between the problem and a model to achieve a correct solution. It does not differentiate well between high level problems, because no matter how many iterations are done it fails to generate the correct code.

In this study, we propose a novel approach to evaluate problem difficulty based on the CoT length generated by reasoning LLMs.
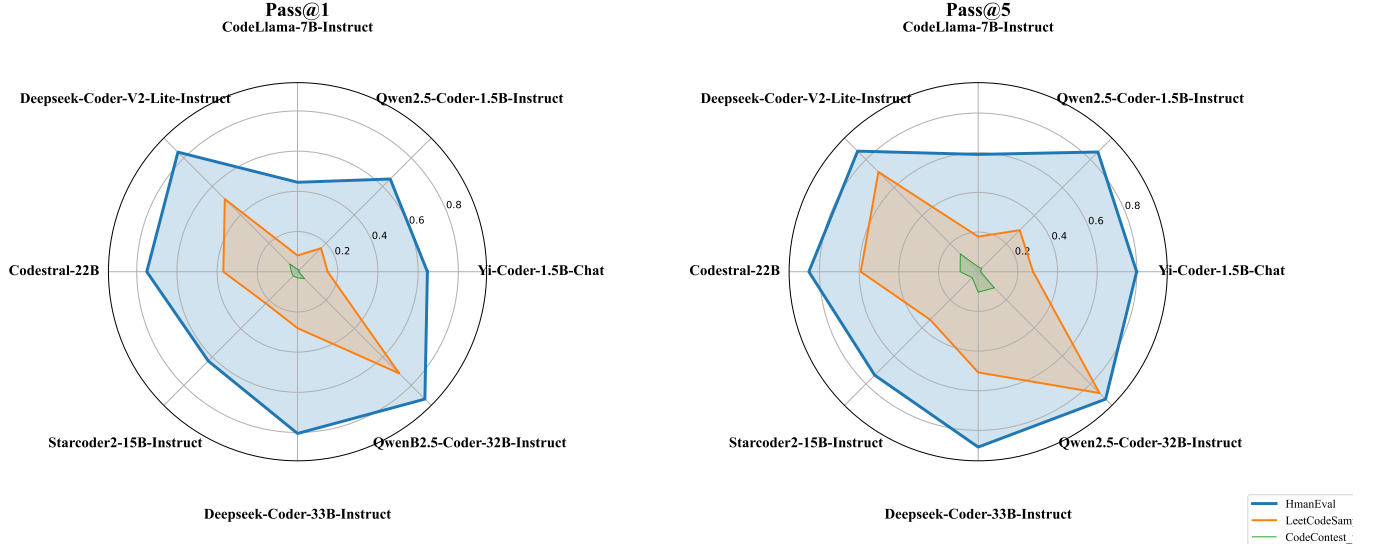
**Figure 2: The performance of eight code LLMs on HumanEval, LeetCodeSample and CodeContests.**

## 2.3 Model Selection

Different LLMs exhibit distinct strengths and weaknesses and no single LLM currently dominates across all tasks. For such issues, the prevalent solution currently is to select the optimal LLM based on the specific task.

Jiang et al. [26] introduced PairRanker and GenFuser, which generate improved outputs by synthesizing results from all LLMs. However, this method requires invoking LLMs $O(T^2)$ times, where $T$ represents the number of models. To optimize both performance and efficiency in LLM selection, researchers have proposed various methods. Cascading strategies [10, 21, 56] sequentially invoke a series of pre-ranked models, typically ordered by capacity from smallest to largest. This process stops when a model's output meets a predefined confidence threshold. Nevertheless, these methods still requires at least $O(T)$ model invocations during inference. Building on these efforts, Chen et al. [12] proposed RouterDC, which employs a routing mechanism to precisely identify the optimal model for a given task. Meanwhile, Bae et al. [6] introduced ComplexityNet, a method that involves fine-tuning an embedded model, DaVinci-002 [7], to select models. It also requires only a single call to the selected model to complete the task.

However, model selection is not only influenced by LLM's capability but also by additional factors such as computational cost and problem difficulty. So, we propose AdaptiveLLM, which integrates multiple factors, including model performance, computational cost, and problem difficulty, into its decision-making process.
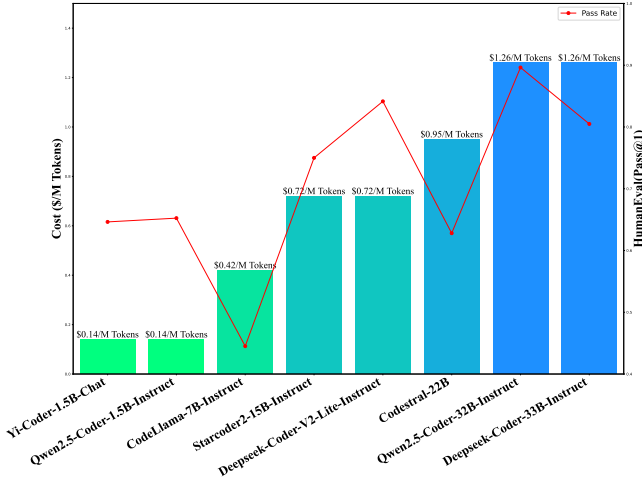
## 3 Preliminary Study

To investigate the capabilities differences among various LLMs and compare their capabilities against pricing, we conduct a preliminary study in this section. We select eight code LLMs: Yi-Coder-1.5B-Chat [54], Qwen2.5-Coder-1.5B-Instruct [24], CodeLlama-7B-Instruct [43], Starcoder2-15B-Instruct [34], DeepSeek-Coder-V2-Lite-Instruct [58], Codestral-22B [47], DeepSeek-Coder-33B-Instruct [20] and Qwen2.5-Coder-32B-Instruct [24]. All these models are

evaluated on three benchmark datasets: HumanEval [11], Leet-CodeSample [33], and CodeContests [29]. Figure 2 illustrates the evaluation results, while Figure 3 shows the comparison of cost ($/M Tokens) and HumanEval pass@1 score across eight models. Our key findings are summarized below.

**Dataset Difficulty Analysis**: The evaluated datasets exhibit a clear gradient of complexity: **HumanEval < LeetCodeSample < CodeContests**. HumanEval focuses on fundamental programming tasks with 164 Python problems that requires function-level code completion based on signature and functional descriptions. Leet-CodeSample, sourced from the LeetCode, presents more complex algorithmic challenges with longer problem descriptions. Code-Contests, derived from Codeforces competitions, which problems involving advanced algorithmic paradigms such as dynamic programming, graph theory, and combinatorial optimization. And These problems have strict time/space constraints and require full-file code generation rather than function completion. Additionally, CodeContests includes challenges with image-based problem descriptions (`<image>` tags), posing additional comprehension barriers for text-only models. As shown in Figure 2, on HumanEval, most models achieve pass@1 score above 60%, while performance drops significantly on CodeContests (average pass@1 < 3%).

**Parameter Size Effects**: We observe a strong positive correlation between model parameter size and code generation performance. For example, increasing the parameters from 1.5B (Yi-Coder) to 32B (Qwen2.5-Coder) yields a 59% improvement in the LeetCode-Sample pass@5 score and 56% improvement in LeetCodeSample pass@1 score. Similarly, on CodeContests, Yi-Coder only achieved 0.61% in pass@1 score and 1.21% in pass@5 score, indicating that it only solved a small number of problems. In contrast, Qwen2.5-Coder obtained 4.85% pass@1 score and 11.52% pass@5 score, reflecting a notable improvement. After conducting a correlation analysis between parameter size and HumanEval pass@1 scores, we found that the correlation coefficient was 0.72, which exceeds 0.7. **This shows that increasing the parameter size brings about**

**significant changes in a model performance.** However, it is worth noting that for models like the 7B CodeLlama, its average performance across the three datasets is significantly lower than that of the two 1.5B models (Yi-Coder-1.5B-Chat and Qwen2.5-Coder-1.5B-Instruct). We attribute this to continuous architectural optimizations. Newer models often have advanced architectures and training techniques, enabling them to outperform older models.



**Figure 3: The relationship between Cost and Performance.**

**Cost-Performance Trade-offs**: Figure 3 reveals a growth in computational overhead as LLM's capacity and parameter size increases. Qwen2.5-Coder-1.5B-Instruct has the highest performance as well as the largest number of parameters, and it also has the highest cost at $1.26/M Tokens. In contrast, CodeLlama-7B-Instruct only need $0.42/M Tokens and it has the worst performance in HumanEval with only 7B parameters, less than Qwen's parameters. Larger models need more memory requirements, with 30B+ models like Qwen2.5-Coder-32B-Instruct needing more than 60GB VRAM for BF16 inference–requiring at least four NVIDIA RTX 4090 GPUs (24GB each). Cloud costs amplify disparities: models with 16.1B+ parameters need costs that are 9× higher than those with 0-4B parameters on Fireworks (https://fireworks.ai/pricing). Although larger models achieve superior code generation performance, their operational costs become prohibitive for budget-constrained projects. **Thus, the cost-performance trade-off underscores the need to select an optimal LLM, balancing parameter scale with financial resources.**

Based on above findings, the design of a framework that can maintain model performance while effectively reducing inference costs has become a critical issue. To address this, we propose a novel framework that aims to leverage the strengths of multiple large models. The core principle of this framework is to dynamically allocate code tasks to the most suitable model based on specific task requirements. For example, simpler tasks are prioritized to smaller parameter-scale models to minimize inference costs, while complex tasks are delegated to larger parameter-scale models to ensure task completion quality. In the following section, we will provide a detailed explanation of the framework.

## 4 Methodology

Figure 4 illustrates the overall architecture of AdaptiveLLM. As a supervised learning approach, AdaptiveLLM first requires reconstructing the three datasets (Section 4.1) due to the absence of learning objectives in their original form. The framework consists of two core components: an embedding layer (Section 4.2) and a classification layer (Section 4.3). Finally, we conduct a comprehensive evaluation of the AdaptiveLLM framework.

### 4.1 Construction of LLM Ranking Datasets

To construct the AdaptiveLLM framework, we first conduct a research study to identify datasets and large models applied to the AdaptiveLLM framework and ranked the models used for each question based on the performance and cost of the model answers. Due to the variety of code LLMs, with their different parameters, architectures, and inference costs, to ensure that the pool of models in the AdaptiveLLM framework can cover most of the LLMs, we select eight representative LLMs among models with different parameter sizes and with different architectures. A brief description of these LLMs is provided in Section 5.1.1. The construction process of LLM ranking dataset is shown in Figure 5. Constructing the LLM ranking dataset requires a total of two steps. The first step is to first have the models in the model candidate pool answer the questions in the selected dataset, and then record metrics such as correctness, token spend, and cost. The second step is to select the optimal LLM for each question based on the recorded metrics, and this step will rank the LLMs used in each question based on the metrics. Finally, the constructed dataset will be divided into training and test sets, and the problem-optimal LLM pairs in the training set will be used to train the XGBOOST classifier in Section 4.3.

*4.1.1 Recording of test result metrics.* In Section 3, we conduct experiments on three datasets for eight large models. To ensure consistency in testing the generated code files, we implement the test validation in SandboxFusion sandbox environment [33] and record the accuracy of each answer for each model. At the same time, we also record the resource consumption metrics of the models, and the price of LLM is provided by SiliconFlow (https://cloud.siliconflow.cn/models), a cloud interfacing platform. However, given that not all of the selected models are deployed on the platform, we use the consumption data of a model with the same parameter size as a proxy for the non-deployed models. Following this step, we obtain a dataset comprising HumanEval, LeetCodeSample, and CodeContests. For each problem in the dataset, it includes five responses from each of the eight LLMs, along with their corresponding correctness results, accuracy rates, and token costs for all five responses.

*4.1.2 LLM ranking.* Since AdaptiveLLM operates as a supervised learning framework, for each problem, we need to predefine an optimal LLM. This optimal LLM should achieve both high accuracy in solving the problem and minimal computational costs. To accomplish this, for every problem, we rank the performance of the eight LLMs on each problem and select the top-ranked model as the optimal LLM. So, we select LLM with the highest $Score_i$ as the optimal model for each problem. $Score_i$ is calculated as follows:
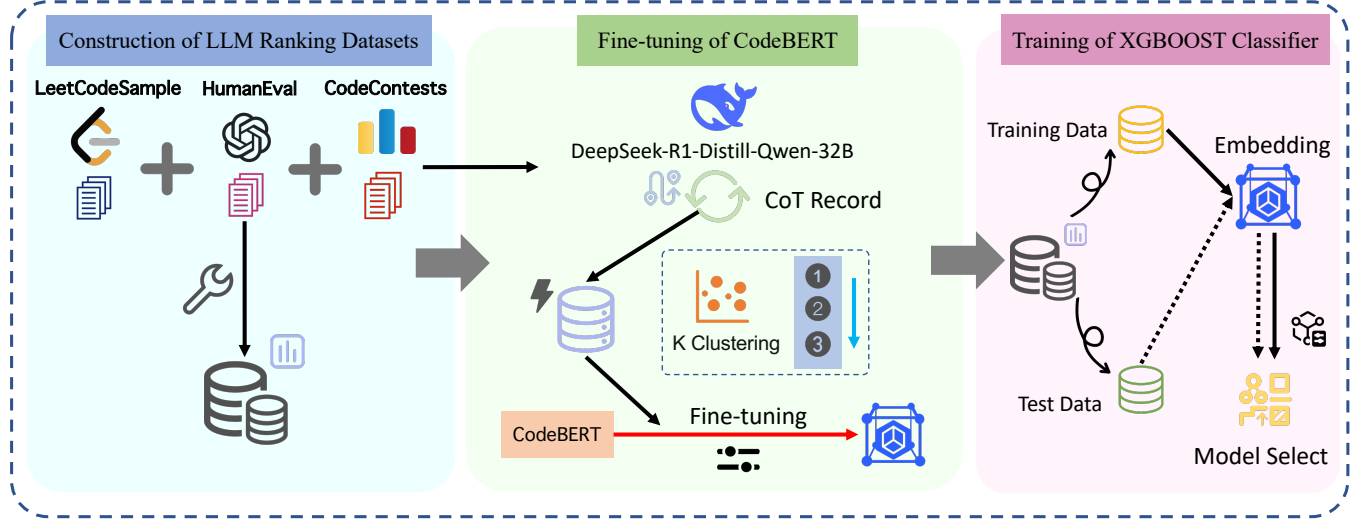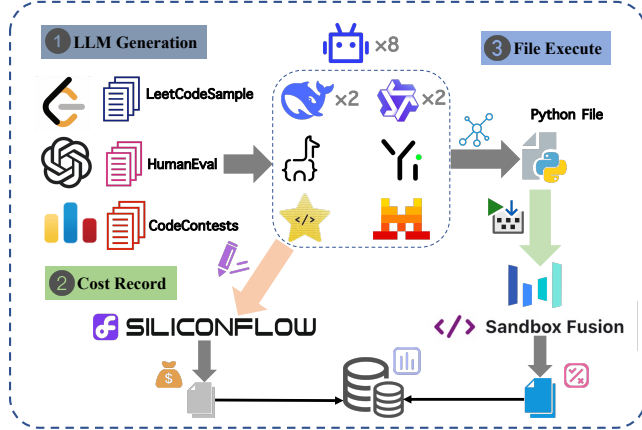
Figure 4: The overall architecture of AdaptiveLLM.



Figure 5: The process of LLM ranking dataset construction.

$$Score_i = \log(\max_{j=1}^{N}(Tokens_j) \times MaxPrice) \times pass_i$$
$$- \log(Tokens_i \times Price_i) \qquad (1)$$

The formula reaches the optimization of model selection with double constraints. First, $pass_i$ is the test pass rate variable, which reflects the proportion of the model that passes all the test cases in five responses. Second, the resource consumption (the product of the number of Token and the inference price) is transformed into a penalty term by introducing a logarithmic operation. This makes the model's score higher the more efficient and less costly its inference is based on the guarantee of correctness of answers. Specially, $MaxPrice$ selects the highest unit price among all models as the normalized benchmark, a move that eliminates the quantitative effects arising from the differences in the pricing systems of different models.

Based on this scoring formula, we can construct a model recommendation system. For each code problem, after calculating $Score_i$

for all candidate models, the model with the highest score is selected as the optimal solution model for the problem. Due to the joint optimization of $Tokens_i$ and $Price_i$ in the formula, the selected model can satisfy the two objectives of code correctness and cost control at the same time. The result dataset covers the complete problem description, model-generated code, token consumption records, inference price details, and optimal model labels based on $Score_i$. You can find detailed information about the dataset in the public repository that we have provided.

## 4.2 Fine-tuning of CodeBERT

In the preliminary study in Section 3, we find that some models perform better when dealing with more complex problems, but at the same time, these models have larger parameter sizes and consume more resources. Given this finding, we would like to incorporate the difficulty of the problem into AdaptiveLLM framework to obtain the best prediction results. However, in real-world scenarios, code problems often lack difficulty annotations. Even when such annotations are available, they may not align with the perception of difficulty as understood by LLMs [28, 42]. To address this challenge, we propose leveraging the length of CoT generated by inference models, such as DeepSeek R1, as a proxy for problem difficulty. This is motivated by the observation that longer CoT lengths generally correspond to higher problem complexity, providing a reliable and automated way to estimate difficulty without relying on manual annotations. So we first use DeepSeek-R1-Distill-Qwen-32B, a distilled version of the reasoning model DeepSeek R1, to assess the difficulty of each problem, and then fine-tune the CodeBERT [16] embedding model using triplet contrast loss.

It is important to note that we do not record cost when using DeepSeek-R1-Distill-Qwen-32B. This is because we only need to evaluate the cost of the model selected in the model candidate pool specifically when answering the question, not the cost of using the model in the training framework phase. Using DeepSeek R1 is so that the question difficulty information is embedded in the

final embedding vector representation. And the involvement of DeepSeek R1 is not required when using AdaptiveLLM.

### 4.2.1 Labeling problem difficulty.
For the constructed dataset, we will use DeepSeek-R1-Distill-Qwen-32B to record the reasoning length for each problem. As seen in Figure 6, problems of different difficulty levels show significant differences in the reasoning length. The reasoning length of a simple problem may be in a lower order of magnitude, while the reasoning length of a complex problem jumps to a higher order of magnitude. For example, the reasoning length of `CodeContests/28` problem is 41 058, that of `LeetCode/22` problem is 24 806, and that of `HumanEval/0` problem is only 6832. This difference reflects the fact that when solving complex problems, the reasoning model needs more complex reasoning steps and logical processes to reach a solution. In contrast, the answer lengths remain largely consistent across problems of varying difficulty. We use DeepSeek-R1-Distill-Qwen-32B and count the average of the reasoning lengths and answer lengths on each dataset. The detailed data is shown in Table 1. This implies that the size of the content output by the model is relatively stable when generating the final code answer, despite the different complexity of the questions.
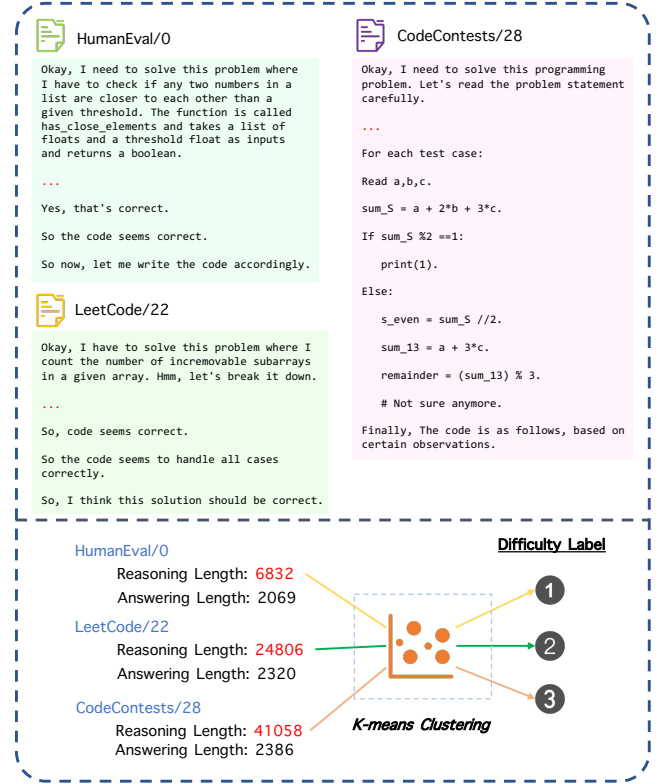
**Table 1: Comparison of DeepSeek-R1-Distill-Qwen-32B's averages of reasoning lengths and answer lengths on different datasets shows a large variability in reasoning lengths, while answer lengths are not discriminatory.**

| Avg Length | HumanEval | LeetCode | CodeContests |
|---|---|---|---|
| Reasoning Len | 7837.46 | 24606.85 | 40882.77 |
| Answering Len | 1985.59 | 2676.79 | 3192.45 |

Considering this, we decided to use the k-means clustering algorithm to cluster and analyze the reasoning lengths of all questions. Following the clustering process, three distinct difficulty levels, labeled as 1, 2, and 3, were identified. That is, we will let the CoT lengths recorded by DeepSeek R1 be used as inputs in all problems, and using the k-means algorithm, we will classify the problems into 3 categories representing three different levels of difficulty: easy, medium, and hard. And then, these classification results will be incorporated into the subsequent CodeBERT fine-tuning process.

### 4.2.2 Contrast loss learning.
The purpose of contrast learning is to allow the CodeBERT embedder to learn the difficulty information of the problem. After contrast learning process, the problem is transformed into an embedding vector that will have difficulty information, which is useful for our subsequent LLM selection.

In contrast learning, optimization process is driven by triplet loss, where each triplet comprises an anchor sample, a positive sample, and a negative sample [23, 50]. Specifically, anchored samples are problems that require optimization of the embedding space, positive samples are semantically of the same difficulty as anchored samples, and negative samples are semantically of a different difficulty than anchored and positive samples. At this point three clusters($C_1, C_2, C_3$) have been obtained after labeling, representing the three different difficulties of the problem. For each anchored problem $x_i$, the following sampling method is used to construct pairs of positive and negative samples:



**Figure 6: The process of CoT difficulty labeling.**

- *Positive sample*: Randomly select a problem with a similar CoT length from the difficulty cluster $C_i$ that is the same as the anchored sample.
- *Negative sample*: Randomly choose one of the two clusters $C_j$ and $C_k$ with different difficulties, and then randomly select a problem with a different CoT length from this cluster.

Then, we can use triplet data to fine-tune the CodeBERT embedding model. The core objective of contrastive learning is to bring queries with similar features and complexities closer together while pushing unrelated queries with different complexities further apart [38]. The design of triplet data enables the model to learn semantic representations by associating the anchor query with its positive counterparts, positioning them closer in the embedding vector space. On the contrary, we expect the model to push unrelated queries farther away from the anchor query.

Specifically, first, we use CodeBERT to extract the embedding representations of each word in every query. These word-level embeddings are aggregated through a pooling layer to generate fixed-size sentence-level embeddings. We denote the mapping of problem $x_i$ into the embedding space $\mathbb{R}^p$ by CodeBERT before fine-tuning as $\mathcal{E}(x_i; \mathbf{w})$. Based on this, we define the following contrastive loss formula:

$$Loss = \max\left(0, \frac{\mathbf{a} \cdot \mathbf{n}}{\|\mathbf{a}\|\|\mathbf{n}\|} - \frac{\mathbf{a} \cdot \mathbf{p}}{\|\mathbf{a}\|\|\mathbf{p}\|} + margin\right)$$
$$= \max\left(0, \mathbf{a} \cdot \mathbf{n} - \mathbf{a} \cdot \mathbf{p} + margin\right) \qquad (2)$$

- $\mathbf{a} = \mathcal{E}(x_{\text{anchor}}; \mathbf{w})$: The embedding of the anchor sample.

- $\mathbf{p} = \mathcal{E}(x_{\text{positive}}; \mathbf{w})$: The embedding of the positive sample.
- $\mathbf{n} = \mathcal{E}(x_{\text{negative}}; \mathbf{w})$: The embedding of the negative sample.

The goal of this loss formula is to learn an embedding space where semantically similar sentences are close to each other, while semantically different sentences are far from each other. The *margin* in the formula is a positive value (by default set to 1), which is used to define the minimum gap between the distance between the anchor and the positive sample and the distance between the anchor and the negative sample. This parameter ensures that negative samples are not just simply pushed away from positive samples but maintain a meaningful distance. In addition, $\max(0, \cdot)$ ensures that the loss value is non-negative. This means that when the distance between the negative sample and the anchor is already large enough, the loss will be zero, indicating that the current triplet does not need further update.

## 4.3  Training of XGBOOST Classifier

Once CodeBERT is fine-tuned, it will serve as a feature extractor for the embedding layer, assisting the XGBoost classifier [13] in supervised learning using the LLM ranking dataset constructed in Section 4.1.

The rationale for selecting XGBoost as the final classifier is evident. Our objective is to construct a framework for selecting the optimal code LLM that ensures both accuracy and minimal resource consumption. Directly incorporating a neural network (NN) into CodeBERT's output layer would contradict our goal due to its substantial computational demands. Therefore, we opted for the lightweight classifier architecture XGBoost, which has been validated as effective in prior research [5, 14].

In our framework, the input features are the embedding vectors generated by the fine-tuned CodeBERT, and the target labels are the optimal model labels corresponding to each problem (specifically, the model that attains the highest score calculated by Score$_i$ in Section 4.1.2). It should be noted that during the training process, the classifier never has access to the test set to avoid data leakage.

## 5  Experimental Setups

### 5.1  Studied Models

*5.1.1  Model candidate pool.* Table 2 shows the code models used in our framework. To ensure that the model pool can cover models with different parameter sizes, we select multiple models ranging from a minimum of 1.5 billion parameters to a maximum of 33 billion parameters. This selection aims to provide suitable code models in response to problems of varying difficulties. The inference prices of the large language models listed in the table are sourced from the SiliconCloud cloud inference platform (https://www.siliconflow.com/en/pricing). However, since some models (such as Starcoder2-15B-Instruct and CodeLlama-7B-Instruct) have not been listed on this platform yet, for these models that are not covered, we assume that their costs are close to those of models with similar parameter sizes. Therefore, we use the prices of models with similar parameter sizes as alternative estimates. For code LLM using, we use RTX 4090 locally to deploy these models. As far as we know, the NVIDIA RTX 4090 GPU offers only 24GB of VRAM, necessitating multi-GPU deployment for models with 15B parameters

or larger. Since inference speed is not a consideration in our evaluation framework and cost metrics are standardized through cloud service provider platforms, the use of single vs. multi-GPU configurations does not affect the experimental outcomes. Considering that there are differences in the default parameter settings of different code models, to ensure the consistency of the experiments, we standardize the model settings: for models with their own default parameters, the parameter settings recommended by the official are adopted; for models without provided default parameters, we uniformly set $do\_sample = True$, $temperature = 0.3$, $top\_p = 0.95$, and $top\_k = 20$. Additionally, all models are implemented using the Hugging Face Transformers library in Python for loading LLMs.

*5.1.2  Reasoning model used in experiments.* Table 2 shows the reasoning models used in our experiments. During the construction of AdaptiveLLM, we employ DeepSeek-R1-Distill-Qwen-32B to measure and record the CoT length. To validate the rationale of using CoT length as a proxy for problem complexity, we conduct experiments in RQ1 in Section 6.1 using multiple DeepSeek R1 variants with different parameter scales. For models with relatively smaller parameter sizes, such as DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B, a single RTX 4090 GPU is deployed for inference. Conversely, models with larger parameter sizes, including DeepSeek-R1-Distill-Qwen-14B and DeepSeek-R1-Distill-Qwen-32B, are used via API, with the API service also provided by SiliconCloud (https://api.siliconflow.cn/v1). Although we utilized both local and cloud service deployment methods for DeepSeek, we assume that their impact on the final CoT length is negligible. This is because the difference in deployment only affects the speed of inferencing, not the quality of the answer.

### 5.2  Metrics and Baselines

*5.2.1  Metrics.*

- **pass@k_score**. This metric measures the probability that at least one correct solution is generated among the top K answers. It is commonly used to evaluate the performance of code generation models. A higher pass@k indicates a better ability of the model to generate correct code within the top K responses.
- **pass@k_token**. This metric represents the average number of tokens consumed per problem during the pass@k test. Tokens are the basic units of input and output in the model, and the quantity of tokens used directly reflects the computational resources required for the test.
- **pass@k_price**. This metric indicates the average cost per problem incurred during the pass@k test and is calculated by multiplying the pass@k_token by the inference price per token. Different models may have different pricing for their tokens, meaning that the same number of tokens could result in significantly different costs depending on the model. By using pass@k_price, we can more accurately reflect the true resource cost of running the pass@k test for each model.

*5.2.2  Baselines.* In our experiments, we choose ComplexityNet [6] as the baseline method for comparison. Its core mechanism involves pre-interactions between LLMs and tasks to determine the complexity of the tasks. Specifically, the framework allows LLMs with varying capabilities (Code Llama, GPT-3.5, and GPT-4) to

**Table 2: Detailed information of LLMs used in our experiments. We do not report the API price for Reasoning Models because these models are used to assess the difficulty of coding problems through CoT length rather than being part of the candidate models for selection.**

| Type | Model | Size | Model Link | API Price($/M tokens) |
|------|-------|------|------------|----------------------|
| Candidate Pool | **Yi-Coder-1.5B** | 1.5B | https://hf.co/01-ai/Yi-Coder-1.5B-Chat | 0.14 |
| | **Qwen2.5-Coder-1.5B-Instruct** | 1.5B | https://hf.co/Qwen/Qwen2.5-Coder-1.5B-Instruct | 0.14 |
| | **CodeLlama-7B-Instruct** | 7B | https://hf.co/meta-llama/CodeLlama-7b-Instruct-hf | 0.42 |
| | **Starcoder2-15B-Instruct** | 15B | https://hf.co/bigcode/starcoder2-15b-instruct-v0.1 | 0.72 |
| | **DeepSeek-Coder-V2-Lite-Instruct** | 16B | https://hf.co/deepseek-ai/DeepSeek-Coder-V2-Lite-instruct | 0.72 |
| | **Codestral-22B** | 22B | https://hf.co/mistralai/Codestral-22B-v0.1 | 0.95 |
| | **DeepSeek-Coder-33B-Instruct** | 33B | https://hf.co/deepseek-ai/deepseek-coder-33b-instruct | 1.26 |
| | **Qwen2.5-Coder-32B-Instruct** | 32B | https://hf.co/Qwen/Qwen2.5-Coder-32B-Instruct | 1.26 |
| Reasoning Models | DeepSeek-R1-Distill-Qwen-1.5B | 1.5B | https://hf.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B | - |
| | DeepSeek-R1-Distill-Qwen-7B | 7B | https://hf.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B | - |
| | DeepSeek-R1-Distill-Qwen-14B | 14B | https://hf.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-14B | - |
| | DeepSeek-R1-Distill-Qwen-32B | 32B | https://hf.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B | - |

attempt solving the tasks multiple times. Based on the correctness of the model outputs, tasks are assigned difficulty labels. These labels are then used by a smaller model, DaVinci-002 [7], to classify and allocate tasks to the most suitable models.

Given the potential unavailability of certain LLMs, we have implemented substitutions in our experimental setup. Specifically, we have employed gpt-4o-2024-11-20 as an alternative for GPT-4 and Qwen2.5-7B-Instruct as an alternative to DaVinci-002. These alternative models are chosen to closely approximate the performance of the models that we used in AdaptiveLLM, thereby ensuring the reliability, fairness and validity of the experimental results.

## 5.3 Implementation Details

In Section 4.2, during the fine-tuning of CodeBERT, we employed the AutoTokenizer and AutoModel classes from the Hugging Face Transformers library. And for encoder-only architectures like BERT, the maximum input sequence length is 512. Consequently, we applied truncation to inputs exceeding this limit, with implementation details documented in our repository. This strategy is rational. Because problem descriptions typically appear at the beginning of input sequences, ensuring that the essential semantic information remains preserved in generated embedding vectors.

In Section 4.2.1, we log the optimal CoT length for reasoning models. Given the variability in the lengths of CoT, we employ the DeepSeek-R1-Distill-Qwen-32B model to generate ten responses for each problem. The median length of these responses is selected as the optimal CoT length for each problem. This approach ensures that the chosen length reflects the complexity of the problem. The same methodology is applied in Section 6.1 to get CoT lengths for reasoning models with varying parameter sizes.

We cap the maximum number of tokens generated by reasoning models at 16 384. This setting allows most problems to generate complete CoT. However, for exceptionally complex problems where the CoT length exceeds this limit, we record the length as 0 and exclude these instances from our dataset. For code models in model candidate pool, we set the maximum number of tokens to 2048, ensuring complete responses for all problems without truncation.

In Section 4.3, during XGBoost training process, the dataset is randomly divided into training and testing subsets at a ratio of

70% to 30%. Specifically, the training set consists of 411 coding problems, while the testing set includes 178 coding problems. For the implementation of XGBoost, we utilized the XGBClassifier class from the xgboost library to train the XGBoost model. Detailed parameter configurations are documented in our repository.

## 6 Experimental Results and Analysis

To assess the effectiveness of AdaptiveLLM, we conduct a thorough evaluation addressing the following research questions:

- **RQ1: Rationality Analysis** - How does the difficulty assessed by the CoT length generated by reasoning models compare with the true difficulty labels? Is it reasonable? Is it affected by the parameter size of the reasoning models?
- **RQ2: Overall Performance** - How does AdaptiveLLM perform on the test set metrics compared to a single LLM and baseline method?
- **RQ3: Ablation Study** - What is the impact of the embedding layer fine-tuning in AdaptiveLLM on the overall performance?
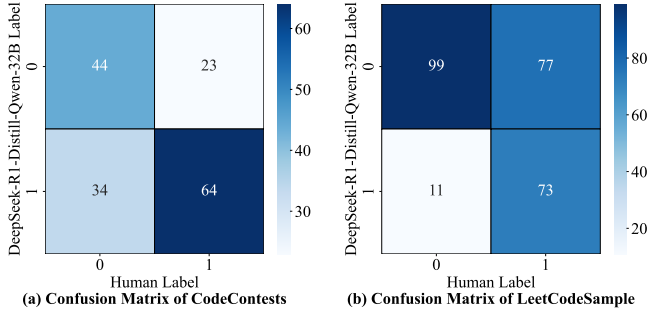
## 6.1 RQ1: Rationality Analysis

*6.1.1 RQ1-1: Is there a significant difference between the difficulty assessed by CoT length and the true difficulty labels?* In the LeetCode-Sample and CodeContests datasets, difficulty labels are assigned based on human performance. Specifically, the difficulty in Leet-CodeSample is determined by the acceptance rate (acRate), and CodeContests relies on competition ratings (cf_rating). These metrics reflect the actual performance of human solvers in terms of success rates.

To enable comparative analysis, we employ the K-means clustering algorithm on each dataset to classify problems with actual difficulty labels into two categories: easy and hard. We then apply the same binary classification to the difficulty labels generated by the DeepSeek-R1-Distill-Qwen-32B model, which are based on the CoT length. Through this process, we compare the model's predicted difficulty against the true difficulty and construct a confusion matrix to analyze.

As shown in Figure 7, even under the binary clustering, there is a significant difference between the difficulty assessed by CoT

**(a) Confusion Matrix of CodeContests**    **(b) Confusion Matrix of LeetCodeSample**
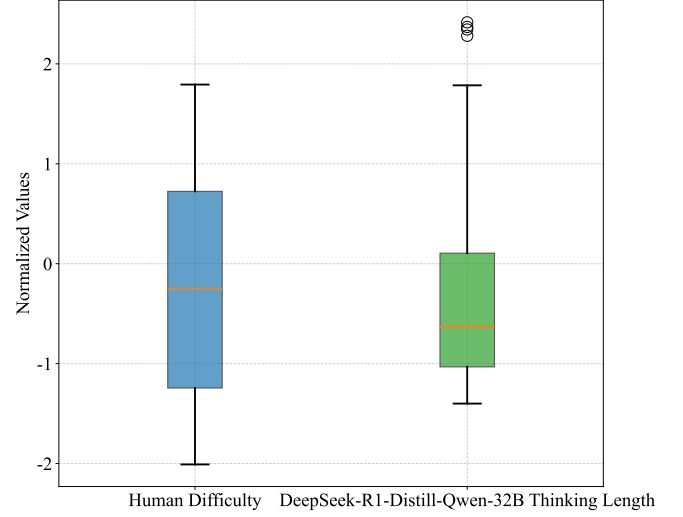
**Figure 7: The confusion matrix between CoT difficulty and true difficulty.**

length and the true difficulty labels. Specifically, in the CodeContests dataset, a total of 57 problems were classified into different difficulty categories by the two methods, while in the LeetCodeSample dataset, this number reached 88. Notably, in the LeetCodeSample dataset, 77 problems labeled as "hard" (label 1) by the true labels were classified as "easy" (label 0) by the CoT length assessment. This finding indicates a substantial difference between the difficulty labels assigned by the R1 model based on CoT length and the original difficulty labels in the dataset.

The original difficulty labels were created by humans, but the actual problem-solving is done by LLMs. If we use the difficulty with which humans view a question to guide LLM in answering the question, it may be biased. So, is LLM's perception of the difficulty of the question is more in line with LLM's perspective of answering the question?

*6.1.2    RQ1-2: Which difficulty assessment method is more practical for LLM?.* Given the inconsistency in difficulty assessment, we design and conduct the following experiment. After logging the CoT lengths, we normalized the extracted CoT lengths using a normal distribution normalization to ensure comparability of the data. Meanwhile, each problem in the LeetCodeSample dataset is accompanied by a difficulty label based on human perception, which is defined as $1 - acRate$. Problems with higher difficulty have lower acceptance rates, resulting in higher difficulty scores. To facilitate comparison with the CoT length, we also normalize these human perception-based difficulty labels using the same normal distribution normalization. Finally, we extract problems from LeetCodeSample that could be correctly answered by Qwen2.5-Coder-32B-Instruct and creat box plot for these problems based on the normalized difficulty data from both methods.

The box plot is shown in Figure 8. Given that the selected data are from problems that the large model can answer correctly, it is expected that the difficulty labels of these problems would generally be low. Meanwhile, the corresponding CoT lengths of the reasoning model would be relatively short. As observed from the box plot, approximately 75% of the problems have CoT lengths less than 0. This indicates that the majority of the difficulty labels based on the reasoning model's CoT lengths are relatively simple, which aligns with our selection of problems that LLMs can correctly answer. However, for the human labeled data, many scores are above 0, indicating that the actual programmers still consider these selected problems to have a certain level of difficulty. This is inconsistent



**Figure 8: The box plot for comparison of CoT difficulty and true difficulty.**

with the performance of LLMs on the questions, as LLMs answer all of these questions correctly.

Moreover, the overall distribution of the difficulty labels based on CoT length is more concentrated compared to those based on human perception. This implies that CoT length tends to reflect that these promblems are generally simple. However, there are some outliers in the CoT length data. These outliers exhibit significantly longer CoT lengths that deviate from the majority of the distribution. This observation underscores that CoT still demonstrates a degree of uncertainty, potentially overcomplicating some otherwise simple problems. But it is still reasonable to apply CoT length for the few cases which can be ignored.

The results overall indicate that using CoT length as a criterion for classifying problem difficulty is rational. More importantly, compared to the human difficulty labels, it may more accurately reflect the actual performance of LLMs in code generation tasks. For example, certain logical reasoning problems that are complex for humans may be relatively straightforward for LLMs, while some challenges that are difficult for LLMs might be easy for humans. Such phenomena are not uncommon in the reasoning processes of LLMs. We believe this finding provides valuable insights and opens new perspectives for future research.

**Table 3: The comparison of CoT length generated by models with different parameter sizes.**

| Model Group | ARI | FMI |
|---|---|---|
| 1.5B - 14B | 0.4650 | 0.6538 |
| 7B - 14B | 0.6540 ↑0.189 | 0.7778 ↑0.124 |
| 1.5B - 32B | 0.3476 | 0.5845 |
| 7B - 32B | 0.4845 ↑0.1369 | 0.6741 ↑0.0896 |
| 14B - 32B | 0.5711 ↑0.2235 | 0.7271 ↑0.1426 |

*6.1.3    RQ1-3: What are the differences in CoT lengths generated by reasoning models of different parameter sizes?* We further investigate whether the size of model parameters affects the CoT length

**Table 4: Performance comparison between AdaptiveLLM and baselines.**

| Model | pass@1 | | | pass@5 | | |
|---|---|---|---|---|---|---|
| | score | token | price($) | score | token | price($) |
| Yi-Coder-1.5B | 26.41% | 329.02 | 4.61e-05 | 35.96% | 1659.47 | 23.23e-05 |
| Qwen2.5-Coder-1.5B-Instruct | 29.78% | 390.71 | 5.47e-05 | 39.33% | 1904.76 | 26.67e-05 |
| CodeLlama-7B-Instruct | 17.98% | 522.68 | 21.95e-05 | 29.21% | 2472.89 | 103.86e-05 |
| Starcoder2-15B-Instruct | 29.21% | 193.47 | 13.93e-05 | 39.89% | 957.97 | 68.97e-05 |
| DeepSeek-Coder-V2-Lite-Instruct | 49.44% | 563.68 | 40.58e-05 | 58.43% | 2837.76 | 204.32e-05 |
| Codestral-22B | 44.38% | 476.48 | 45.27e-05 | 52.81% | 2383.47 | 226.43e-05 |
| DeepSeek-Coder-33B-Instruct | 39.89% | 359.61 | 45.31e-05 | 49.44% | 1839.34 | 231.76e-05 |
| Qwen2.5-Coder-32B-Instruct | 57.87% | 763.57 | 96.21e-05 | 71.35% | 3834.90 | 483.20e-05 |
| GPT3.5 | 35.39% | 130.42 | 19.56e-05 | 45.51% | 646.28 | 96.94e-05 |
| GPT4o | 60.67% | 551.62 | 551.62e-05 | 72.47% | 2741.79 | 2741.79e-05 |
| ComplexityNet | 37.08% | 380.30 | 256.83e-05 | 50.56% | 1951.39 | 1306.66e-05 |
| **AdaptiveLLM** | **44.94%**↑7.86 | **428.80**↑48.5 | **28.49e-05**↓228.34 | **56.18%**↑5.62 | **2094.98**↑143.59 | **140.07e-05**↓1166.59 |
| AdaptiveLLM (w/o finetune) | 43.82%↓1.12 | 411.84↓16.96 | 26.74e-05↓1.75 | 55.62%↓0.56 | 2037.09↓57.89 | 133.99e-05↓6.08 |

and its corresponding difficulty classification. To address this question, we select DeepSeek-R1-Distill-Qwen models with varying parameter sizes (1.5B, 7B, 14B, and 32B), extract the CoT lengths generated by each model, and conduct clustering analysis on these data. Subsequently, we compare the difficulty classifications of different parameter size models in a pairwise manner.

The results reveal that the consistency of difficulty classification significantly increases as the model parameter size grows. In the clustering comparison among the 1.5B, 7B, 14B, and 32B reasoning models, it is observed that as the parameter size increases, their classification results gradually converge with those of the 32B R1 reasoning model. This trend is quantified using the Adjusted Rand Index (ARI) and Fowlkes-Mallows Index (FMI), as shown in Table 3. These metrics measure the similarity between the clustering results of models with different parameter sizes. Results show that model group with closer parameter sizes achieve higher ARI and FMI values, indicating stronger consistency in difficulty classification. For example, the ARI and FMI values between the 14B and 32B models are 0.5711 and 0.7271, respectively, while the corresponding values between the 1.5B and 32B models are only 0.3476 and 0.5845.

This finding suggests that smaller models exhibit randomness in CoT length, leading to less stable classifications. In contrast, larger models show a stronger correlation between CoT length and actual difficulty, improving classification stability and reliability.

> *Answer to RQ1:* The difficulty assessment based on CoT length of reasoning models differs from human. Nevertheless, the CoT length method seems to better reflect the actual performance of LLMs when solving problems. And reasoning models with larger parameter sizes tend to generate more stable and consistent CoT.

## 6.2 RQ2: Overall Performance

*6.2.1 RQ2-1: How does the overall performance of AdativeLLM compare with single LLM?.* Table 4 presents a detailed comparison of key metrics, including pass@1 and pass@5 score, token consumption, and inference cost, for eight code LLMs and the AdaptiveLLM framework on the benchmark dataset.

Results demonstrate that AdaptiveLLM framework exhibits significant performance advantages. In terms of accuracy, the AdaptiveLLM framework achieves pass@1 score of 44.94% and pass@5 score of 56.18%. These results not only surpass those of the Codestral-22B model with 22B parameters (pass@1 score of 44.38% and pass@5 score of 55.28%), but also show a clear superiority compared to the larger-parameter DeepSeek-Coder-33B-Instruct model (pass@1 score of 39.89% and pass@5 score of 49.44%). Notably, while maintaining a high level of accuracy, the token consumption of the AdaptiveLLM framework (pass@1: 428.80 and pass@5: 2094.98) is significantly lower than that of most larger-parameter models, such as Qwen2.5-Coder-32B-Instruct (pass@1: 763.57 and pass@5: 3834.90).

In terms of cost-effectiveness, AdaptiveLLM framework demonstrates remarkable economic efficiency. Its inference cost (pass@1: 28.49e-05$ and pass@5: 140.07e-05$) is comparable to CodeLlama-7B-Instruct model (pass@1: 21.95e-05$ and pass@5: 103.86e-05$), while achieving significantly higher accuracy. Compared to larger-parameter models, the cost control advantage of the AdaptiveLLM framework is even more evident. For example, compared to the Codestral-22B model (pass@1: 45.27e-05$ and pass@5: 226.43e-05$), the AdaptiveLLM framework reduces inference cost by approximately 37% while maintaining a certain level of accuracy.

*6.2.2 RQ2-2: How does the overall performance of AdativeLLM compare with sota baseline method?* To comprehensively evaluate AdaptiveLLM against the state-of-the-art baseline method ComplexityNet [6], we analyze both accuracy and cost-efficiency metrics. AdaptiveLLM achieves a pass@1 score of 44.94%, outperforming ComplexityNet by 7.86%, while consuming moderately more tokens (428.80 vs. 380.30). This accuracy gain is maintained in pass@5, where AdaptiveLLM scores 56.18%, surpassing ComplexityNet by 5.62%. The token increase (143.59 tokens for pass@5) is strategically offset by the framework's dynamic model selection, which prioritizes cost-effective inference without compromising performance.

The most striking advantage lies in cost reduction. For pass@1, AdaptiveLLM reduces inference costs by 88.3% compared to ComplexityNet, achieved through adaptive routing of tasks to smaller

LLMs (e.g., Yi-Coder-1.5B) for simple cases and reserving larger models (e.g., Qwen2.5-Coder-32B) for complex scenarios. In pass@5, the cost reduction reaches 88.9%.

Further analysis reveals that AdaptiveLLM achieves 82.5% cost reduction compared to GPT-4o (28.49e-05$ vs. 551.62e-05$ in pass@1) while retaining 73.6% of GPT-4o's accuracy (44.94% vs. 60.67%). This explains why ComplexityNet incurs significantly higher resource consumption — it relies heavily on closed-source models like GPT, whose high API pricing drives up inference costs compared to AdaptiveLLM. While closed-source models demonstrate better accuracy performance compared with high-parameter open-source models, their inference costs are substantially higher. This difference highlights a critical insight for future model selection: prioritizing open-source models in candidate pools can achieve competitive performance at a fraction of the cost, offering a more sustainable approach to balancing accuracy and economic efficiency.

> *Answer to RQ2:* AdaptiveLLM achieves performance with 44.94% pass@1 and 56.18% pass@5 scores, outperforming both single large LLM and ComplexityNet baseline by +7.86% in pass@1. It significantly reduces inference costs compared to ComplexityNet, demonstrating exceptional cost-effectiveness.

## 6.3 RQ3: Ablation Study

Through the design of ablation experiments, we compared and analyzed the performance differences of the AdaptiveLLM framework in terms of pass@1 and pass@5 accuracy, token consumption, and inference cost before and after CodeBERT fine-tuning. As shown in Table 4, the AdaptiveLLM framework fine-tuned with difficulty-aware tuning exhibited significant improvements across all metrics.

In terms of accuracy, the fine-tuned framework achieved a 1.12 percentage point increase in the pass@1 metric (from 43.82% to 44.94%) and a 0.56 percentage point increase in the pass@5 metric (from 55.62% to 56.18%). This improvement indicates that difficulty-aware fine-tuning effectively enhanced the framework's ability to understand the complexity of problems, thereby enabling it to more accurately select appropriate models for handling problems of varying difficulty levels.

Regarding computational efficiency and cost-effectiveness, the fine-tuned framework demonstrated superior token consumption and inference cost efficiency. Although there was a slight increase in token consumption (pass@1 from 411.84 to 428.80, pass@5 from 2037.09 to 2094.98) and a marginal rise in inference cost (pass@1 from 26.74e-05 to 28.49e-05, pass@5 from 133.99e-05 to 140.07e-05), this increase was positively correlated with the improvement in accuracy, indicating that the additional token and inference cost expenditure brought about effective performance gains. Furthermore, in-depth analysis revealed that for more difficult problems, the AdaptiveLLM framework without fine-tuning still tended to select smaller-parameter models such as Yi-Coder-1.5B-Chat and Qwen2.5-Coder-1.5B-Instruct, which may have led to insufficient capability in handling complex problems. In contrast, the framework fine-tuned with difficulty-aware tuning was able to more accurately identify problem difficulty and accordingly select larger-parameter models suitable for complex problems, thereby significantly enhancing overall performance.

> *Answer to RQ3:* The experimental results confirmed the effectiveness of the fine-tuning strategy based on CodeBERT. By incorporating problem difficulty label information, the framework was able to better understand task characteristics and make more optimal model selection decisions.

## 7 Threats to Validity

**Internal validity.** One potential threat is the randomness in the generation of CoT by reasoning LLMs. To mitigate this, we generated ten responses for each problem and used the median length of these responses as the CoT length for that problem. Another concern is that, due to resource constraints, for a small number of extremely complex problems, the CoT length exceeded the maximum token limit set for model generation. However, since such samples were rare, they were excluded from the dataset. While this exclusion may introduce some bias, its impact is minimal given the small proportion of affected samples.

**External validity.** In our CoT difficulty assessment method, the datasets we selected are primarily focused on single-file code generation tasks and do not involve multi-file interactions or collaborations. As a result, our method may face challenges when applied to project-level code generation tasks.

**Construct validity.** In our experiments, model costs were sourced from the SiliconFlow cloud platform. For models not deployed on this platform, we approximated their costs using models with similar parameter sizes. However, inference costs depend not only on parameter size but also on architecture and inference strategies. For instance, under the same parameter size, Mixture-of-Experts (MoE) models [30, 31] typically have lower inference costs due to their reduced number of active parameters during runtime. Despite these limitations, parameter size is still a practical proxy for cost estimation, acknowledging the potential for minor bias.

## 8 Conclusion

In this paper, we propose AdaptiveLLM, a framework that uses the CoT length of reasoning models to evaluate problem complexity and dynamically select code LLMs. To construct the dataset, we integrate samples of varying difficulty levels and employ a scoring formula to identify the optimal LLM for each problem as ground-truth annotations. For complexity labeling, we utilize the reasoning model DeepSeek-R1-Distill-Qwen-32B to generate CoT sequences, record their lengths, and apply k-means clustering to categorize problems into three difficulty tiers. Based on these difficulty labels, we perform triplet contrastive fine-tuning on the embedding layer of CodeBERT, enabling problem embeddings to encode complexity-aware features that enhance model selection. Finally, we split the constructed dataset into training and test sets to train an XGBoost classifier. Experimental results validate the rationality of difficulty assessment method based on CoT length and effectiveness of AdaptiveLLM framework. We release all code publicly at https://github.com/cjhCoder7/AdaptiveLLM.

## Acknowledgments

# 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] Anysphere. 2009. Codeforces. https://codeforces.com/

[3] Anysphere. 2015. LeetCode. https://leetcode.com/

[4] Anysphere. 2023. Cursor. https://www.cursor.com/

[5] Zeliha Ergul Aydin and Zehra Kamisli Ozturk. 2021. Performance analysis of XGBoost classifier with missing data. *Manchester Journal of Artificial Intelligence and Applied Sciences (MJAIAS)* 2, 02 (2021), 2021.

[6] Henry Bae, Aghyad Deeb, Alex Fleury, and Kehang Zhu. 2023. Complexitynet: Increasing llm inference efficiency by learning task complexity. *arXiv preprint arXiv:2312.11511* (2023).

[7] 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.

[8] G Ann Campbell. 2018. Cognitive Complexity-A new way of measuring understandability. *SonarSource SA* 10 (2018).

[9] Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. 2022. Multipl-e: A scalable and extensible approach to benchmarking neural code generation. *arXiv preprint arXiv:2208.08227* (2022).

[10] Lingjiao Chen, Matei Zaharia, and James Zou. 2023. Frugalgpt: How to use large language models while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176* (2023).

[11] 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).

[12] Shuhao Chen, Weisen Jiang, Baijiong Lin, James Kwok, and Yu Zhang. 2024. RouterDC: Query-based router by dual contrastive learning for assembling large language models. *Advances in Neural Information Processing Systems* 37 (2024), 66305–66328.

[13] Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. 785–794.

[14] Zhuo Chen, Fu Jiang, Yijun Cheng, Xin Gu, Weirong Liu, and Jun Peng. 2018. XGBoost classifier for DDoS attack detection and analysis in SDN-based cloud. In *2018 IEEE international conference on big data and smart computing (bigcomp)*. IEEE, 251–256.

[15] Christof Ebert, James Cain, Giuliano Antoniol, Steve Counsell, and Phillip Laplante. 2016. Cyclomatic complexity. *IEEE software* 33, 6 (2016), 27–29.

[16] 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).

[17] Github. 2021. Github Copilot. https://github.com/features/copilot

[18] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

[19] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[20] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming–The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196* (2024).

[21] Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkrittum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. 2024. Language model cascades: Token-level uncertainty and beyond. *arXiv preprint arXiv:2404.10136* (2024).

[22] T Hariprasad, G Vidhyagaran, K Seenu, and Chandrasegar Thirumalai. 2017. Software complexity analysis using halstead metrics. In *2017 international conference on trends in electronics and informatics (ICEI)*. IEEE, 1109–1113.

[23] Elad Hoffer and Nir Ailon. 2015. Deep metric learning using triplet network. In *Similarity-based pattern recognition: third international workshop, SIMBAD 2015, Copenhagen, Denmark, October 12-14, 2015. Proceedings 3*. Springer, 84–92.

[24] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186* (2024).

[25] Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C Park. 2024. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. *arXiv preprint arXiv:2403.14403* (2024).

[26] Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. *arXiv preprint arXiv:2306.02561* (2023).

[27] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770* (2023).

[28] Bonan Kou, Shengmai Chen, Zhijie Wang, Lei Ma, and Tianyi Zhang. 2023. Is model attention aligned with human attention? an empirical study on large language models for code generation. *arXiv preprint arXiv:2306.01220* (2023).

[29] 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.

[30] Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, et al. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434* (2024).

[31] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

[32] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2023), 21558–21572.

[33] Siyao Liu, He Zhu, Jerry Liu, Shulin Xin, Aoyan Li, Rui Long, Li Chen, Jack Yang, Jinxiang Xia, ZY Peng, et al. 2024. Fullstack bench: Evaluating llms as full stack coder. *arXiv preprint arXiv:2412.00535* (2024).

[34] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173* (2024).

[35] Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583* (2023).

[36] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511* (2022).

[37] Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. 2023. Octopack: Instruction tuning code large language models. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.

[38] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).

[39] OpenAI. 2024. OpenAI o1. https://openai.com/o1/

[40] OpenAI. 2025. OpenAI o3mini. https://openai.com/index/openai-o3-mini/

[41] Gabriel Orlanski, Kefan Xiao, Xavier Garcia, Jeffrey Hui, Joshua Howland, Jonathan Malmaud, Jacob Austin, Rishabh Singh, and Michele Catasta. 2023. Measuring the impact of programming language distribution. In *International Conference on Machine Learning*. PMLR, 26619–26645.

[42] Shuyin Ouyang, Jie M Zhang, Mark Harman, and Meng Wang. 2025. An empirical study on the non-determinism of chatgpt in code generation. *ACM Transactions on Software Engineering and Methodology* 34, 2 (2025), 1–28.

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

[44] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530* (2024).

[45] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295* (2024).

[46] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. 2025. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599* (2025).

[47] Mistral AI team. 2024. Codestral. https://mistral.ai/news/codestral

[48] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).

[49] Chung-Yu Wang, Alireza DaghighFarsoodeh, and Hung Viet Pham. 2024. Selection of Prompt Engineering Techniques for Code Generation through Predicting Code Complexity. *arXiv preprint arXiv:2409.16416* (2024).

[50] Moshi Wei, Nima Shiri Harzevili, Yuchao Huang, Junjie Wang, and Song Wang. 2022. Clear: contrastive learning for api recommendation. In *Proceedings of the 44th International Conference on Software Engineering.* 376–387.

[51] Kurt D Welker. 2001. The software maintainability index revisited. *CrossTalk* 14 (2001), 18–21.

[52] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115* (2024).

[53] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122* (2024).

[54] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652* (2024).

[55] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284* (2023).

[56] Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu Yao. 2023. Large language model cascades with mixture of thoughts representations for cost-efficient reasoning. *arXiv preprint arXiv:2310.03094* (2023).

[57] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual benchmarking on humaneval-x. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining.* 5673–5684.

[58] Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. 2024. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931* (2024).