

# EquiBench: Benchmarking Code Reasoning Capabilities of Large Language Models via Equivalence Checking

Anonymous ACL submission

## Abstract

Equivalence checking, i.e., determining whether two programs produce identical outputs for all possible inputs, underpins a broad range of applications, including software refactoring, testing, and optimization. We present the task of equivalence checking as a new way to evaluate the code reasoning abilities of large language models (LLMs). We introduce EquiBench, a dataset of 2400 program pairs spanning four programming languages and six equivalence categories. These pairs are systematically generated through program analysis, compiler scheduling, and superoptimization, covering nontrivial structural transformations that demand deep semantic reasoning beyond simple syntactic variations. Our evaluation of 17 state-of-the-art LLMs shows that OpenAI o3-mini achieves the highest overall accuracy of 78.0%. In the most challenging categories, the best accuracies are 62.3% and 68.8%, only modestly above the 50% random baseline for binary classification, indicating significant room for improvement in current models' code reasoning capabilities.

## 1 Introduction

Programming has emerged as a key application domain for large language models (LLMs), enabling tasks such as program synthesis (Chen et al., 2021; Austin et al., 2021; Jain et al., 2024), test generation (Yang et al., 2024a), bug detection (Yang et al., 2023), program repair (Xia et al., 2023), and code optimization (Shypula et al., 2023). Recently, there has been growing interest in evaluating how well LLMs can reason about the semantics of code (Ni et al., 2024; Liu et al., 2023; Gu et al., 2024; Chen et al., 2024a; Liu et al., 2024b), i.e., predicting program properties without running the program.

This paper introduces the task of **equivalence checking** as a new way to evaluate the code reasoning capabilities of LLMs. A classic challenge

```
int main() { | int main() { | int main() {
  int a, b;   | int a, b;   | int a, b;
  ... if (a < b) { | ... if (b > a) { | ... if (b <= a) {
    a = 1;    | a = 1;    | a = 1;
  }          | }          | }
  ...      | ...      | ...
}          | }          | }
```

Equivalent                      Inequivalent

Figure 1: **An equivalent and an inequivalent program pair constructed using prior techniques.** Prior works generate such pairs through *basic statement-level syntactic modifications* with minimal semantic reasoning, whereas our approach, presented later, relies on structural program transformations that require much deeper semantic reasoning.

in programming languages and verification, equivalence checking involves determining whether two programs produce identical outputs for all possible inputs. Figure 1 presents examples of equivalent and inequivalent program pairs.

Compared to prior code reasoning tasks, evaluating LLMs using equivalence checking offers distinct advantages. Most notably, it presents a significantly more challenging benchmark than previous tasks, enabling a more rigorous assessment of LLMs' code reasoning capabilities. Equivalence checking requires LLMs to reason over *all possible inputs*, while prior work often focuses on *a single input*, such as output prediction, input prediction (Gu et al., 2024), input-specific program state prediction and execution simulation (Liu et al., 2023; Chen et al., 2024a; Ding et al., 2024; La Malfa et al., 2024; Ni et al., 2024).

Moreover, equivalence checking underpins a broad range of downstream applications, including software refactoring (Pailoor et al., 2024), software testing (Tian et al., 2024), and program optimization (Shypula et al., 2021), surpassing the scope of prior reasoning tasks. By requiring a deep understanding of program semantics and reasoning over all possible inputs, equivalence checking en-

ables the analysis of an expressive range of program behaviors, even including many undecidable problems. Therefore, LLMs that perform well on equivalence checking are likely to be well-suited for tackling more complex programming tasks.

Our proposal requires a benchmark consisting of both equivalent and inequivalent program pairs covering different aspects of equivalence reasoning with varying degrees of difficulty. A large benchmark is essential, making it desirable to automate the benchmark generation process. Existing methods (Badihi et al., 2021; Maveli et al., 2024) mostly rely on *local syntactic changes* such as operand swaps (e.g., changing  $a < b$  to  $b > a$  for equivalent pairs or  $b \leq a$  for inequivalent pairs; see Figure 1), which do not require deep semantic reasoning. However, these approaches are insufficient for benchmarking the equivalence reasoning capabilities of state-of-the-art LLMs. As many existing benchmarks have become saturated (Phan et al., 2025), a more challenging dataset is needed to rigorously assess LLMs’ semantic reasoning abilities.

In this work, we introduce **EquiBench**, a new dataset of 2400 program pairs for equivalence reasoning. EquiBench spans four programming languages—Python, C, CUDA, and x86-64 assembly—providing a systematic benchmark to evaluate LLMs’ code reasoning abilities.

The key technical challenge is to automatically generate (in)equivalent program pairs that demand *deep semantic reasoning beyond simple syntactic variations*. We propose several techniques to achieve this. First, to confirm that basic syntactic variations are well within the reasoning capabilities of state-of-the-art LLMs, we construct an equivalence category based on variable renaming, which barely requires semantic reasoning. Next, we generate equivalent programs by removing dead code, leveraging program analysis to go beyond trivial syntactic changes. By incorporating alias analysis and path feasibility analysis, we increase the difficulty of semantic reasoning in an automated manner. For GPU programs written in CUDA, we generate equivalent pairs by exploring different compiler scheduling strategies, such as loop tiling and shared memory caching, which involve structural transformations that extend far beyond statement-level modifications. We also use *superoptimization* to explore optimal instruction sequences beyond standard compiler optimizations, enabling more aggressive code restructuring. Finally, we include pairs with different algorithmic choices using sub-

missions from online programming platforms.

Our experiments show that EquiBench is a challenging benchmark for LLM-based equivalence checking. Among the 17 models evaluated, OpenAI o3-mini performs best overall, yet achieves only 59.0% in the CUDA category despite achieving the highest overall accuracy of 78.0%. For the two most difficult categories, the best accuracy across all models is 62.3% and 68.8%, respectively. These numbers are only *modestly better than the random baseline*—i.e., 50% accuracy for binary classification. Further analysis shows that variable renaming, a purely syntactic modification, is the easiest equivalence category for models, with accuracy as high as 91.2%. We also find that models are *biased* toward classifying programs with significant structural, non-local transformations as inequivalent. Moreover, prompting strategies such as few-shot in-context learning and Chain-of-Thought (CoT) prompting *barely* enhance LLMs’ semantic reasoning capabilities in equivalence checking, underscoring the fundamental difficulty of the task.

In summary, our contributions are as follows:

- **New Task and Dataset:** We introduce equivalence checking as a new task to assess LLMs’ code reasoning capabilities. We present *EquiBench*, a benchmark for semantic equivalence checking spanning four languages and six equivalence categories.
- **Automated Generation:** We develop a fully automated pipeline to construct diverse (in)equivalent program pairs, using techniques from program analysis, compiler scheduling, and superoptimization. The pipeline covers transformations including syntactic changes, structural modifications, and algorithmic equivalence.
- **Evaluation and Analysis:** We evaluate 17 state-of-the-art models on EquiBench, with the highest overall accuracy reaching 78.0%. In the two most challenging categories, the best accuracy across all models is 62.3% and 68.8%, indicating significant room for improvement. Additionally, we analyze performance across different equivalence categories and prompting strategies.

## 2 Related Work

**LLM Reasoning** Extensive research has evaluated LLMs’ reasoning capabilities across diverse

tasks (Cobbe et al., 2021; Huang and Chang, 2022; Bubeck et al., 2023; Mirzadeh et al., 2024; Zhou et al., 2022; Ho et al., 2022; Wei et al., 2022; Chen et al., 2024b; Clark et al., 2018; Zhang et al., 2024). In the context of code reasoning, i.e., predicting a program’s execution behavior without running it, CRUXEval (Gu et al., 2024) focuses on input-output prediction, while CodeMind (Liu et al., 2024b) extends evaluation to natural language specifications. Another line of work seeks to improve LLMs’ code simulation abilities through prompting (La Malfa et al., 2024) or targeted training (Liu et al., 2023; Ni et al., 2024; Ding et al., 2024). Unlike prior work that evaluates LLMs on predicting program behavior for a specific input, our new code reasoning task and benchmark for equivalence checking assesses LLMs’ ability to reason about all possible inputs.

**Equivalence Checking** Equivalence checking underpins applications such as performance optimization (Shypula et al., 2023; Cummins et al., 2023, 2024), code transpilation (Lu et al., 2021; Yang et al., 2024b; Ibrahimzada et al., 2024; Pan et al., 2024), refactoring (Pailoor et al., 2024), and testing (Felsing et al., 2014; Tian et al., 2024). Due to its undecidable nature, no algorithm can decide program equivalence for all program pairs while always terminating. Existing techniques (Sharma et al., 2013; Dahiya and Bansal, 2017; Gupta et al., 2018; Mora et al., 2018; Churchill et al., 2019; Badihi et al., 2020) focus on specific domains, such as SQL query equivalence (Zhao et al., 2023; Ding et al., 2023; Singh and Bedathur, 2024). EQBENCH (Badihi et al., 2021) and SeqCoBench (Maveli et al., 2024) are the main datasets for equivalence checking but have limitations. EQBENCH is too small (272 pairs) for LLM evaluation, while SeqCoBench relies only on statement-level syntactic changes (e.g., renaming variables). In contrast, our work introduces a broader set of equivalence categories and structural transformations, creating a more systematic and challenging benchmark for assessing LLMs’ semantic reasoning capabilities.

### 3 Benchmark Construction

While we have so far discussed only the standard notion of equivalence (that two programs produce the same output on any input), there are other, more precise definitions of equivalence used for each category in the benchmark. For each category, we

<pre>char b[2]; static int c = 0;  int main() {     char* p1 = &amp;b[0];     int* p2 = &amp;c;     ...     if (p1 == p2) {         c = 1; //dead code     }     ...     return 0; }</pre>	<pre>char b[2]; static int c = 0;  int main() {     char* p1 = &amp;b[0];     int* p2 = &amp;c;     ...     if (true) {         c = 1; //live code     }     ...     return 0; }</pre>
--	--

Figure 2: **An inequivalent pair from the DCE category in EquiBench.** In the left program, `c = 1` is dead code and has no effect on the program state, whereas in the right program, it is executed and alters the program state. Such cases are generated using the Dead Code Elimination (DCE) pass in compilers.

provide the definition of equivalence, which is included in the prompt when testing LLM reasoning capabilities. We describe the process of generating (in)equivalent pairs for the following six categories:

- **DCE:** C program pairs generated via the compiler’s dead code elimination (DCE) pass (Section 3.1).
- **CUDA:** CUDA program pairs created by applying different scheduling strategies using a tensor compiler (Section 3.2).
- **x86-64:** x86-64 assembly program pairs generated by a superoptimizer (Section 3.3).
- **OJ\_A, OJ\_V, OJ\_VA:** Python program pairs from online judge submissions, featuring algorithmic differences (OJ\_A), variable-renaming transformations (OJ\_V), and combinations of both (OJ\_VA) (Section 3.4).

#### 3.1 Pairs from Program Analysis (DCE)

Dead code elimination (DCE), a compiler pass, removes useless program statements. After DCE, remaining statements in the modified program naturally *correspond* to those in the original program.

**Definition of Equivalence.** Two programs are considered equivalent if, when executed on the same input, they *always* have identical *program states* at all corresponding points reachable by program execution. We expect language models to identify differences between the two programs, align their states, and determine whether these states are consistently identical.

```

__global__ void GEMV(const float* A,
                    const float* x,
                    float* y,
                    int R,
                    int C) {
    // Calculate the row index
    // assigned to the thread
    int r = blockIdx.x * blockDim.x
        + threadIdx.x;

    // Return if out of bounds
    if (r >= R) return;
    float s = 0.0f;

    for (int c = 0; c < C; c++) {
        s += A[r * C + c] * x[c];
    }

    y[r] = s;
}

__global__ void GEMV(const float* A, const float* x,
                    float* y, int R, int C) {
    __shared__ float tile[32]; // tiling with shared memory
    int r = blockIdx.x * blockDim.x + threadIdx.x;
    bool valid = (r < R);
    float s = 0.0f;
    for (int start = 0; start < C; start += 32) {
        for (int i = threadIdx.x; i < 32; i += blockDim.x) {
            int c = start + i;
            if (c < C) tile[i] = x[c]; // load x into tile
        }
        __syncthreads();
        if (valid) {
            for (int j = 0; j < min(32, C - start); j++) {
                s += A[r * C + (start + j)] * tile[j];
            }
        }
        __syncthreads();
    }
    if (valid) y[r] = s;
}

```

Figure 3: **An equivalent pair from the CUDA category in EquiBench.** Both programs perform matrix-vector multiplication ( $y = Ax$ ). The right-hand program uses *shared memory tiling* to improve performance. Tensor compilers are utilized to explore different *scheduling strategies*, automating the generation.

**Example.** Figure 2 illustrates an inequivalent pair of C programs. In the left program, the condition ( $p1 == p2$ ) compares the memory address of the first element of the array  $b$  with that of the static variable  $c$ . Since  $b$  and  $c$  reside in different memory locations, this condition can never be satisfied. As a result, the assignment  $c = 1$  is never executed in the left program but is executed in the right program. This difference in program state during execution renders the pair inequivalent.

**Automation.** This reasoning process is automated by compilers through *alias analysis*, which statically determines whether two pointers can reference the same memory location. Based on this analysis, the compiler’s *Dead Code Elimination (DCE)* pass removes code that does not affect program semantics to improve performance.

**Dataset Generation.** We utilize CSmith (Yang et al., 2011) to create an initial pool of random C programs. Building on techniques from prior compiler testing research (Theodoridis et al., 2022), we implement an LLVM-based tool (Lattner and Adve, 2004) to classify code snippets as either dead or live. Live code is further confirmed by executing random inputs with observable side effects. Equivalent program pairs are generated by eliminating dead code, while inequivalent pairs are generated by removing live code.

### 3.2 Pairs from Compiler Scheduling (CUDA)

**Definition of Equivalence.** Two CUDA programs are considered equivalent if they produce the same mathematical output for any valid input, *disregarding floating-point rounding errors*. This definition *differs* from that in Section 3.1, as it does not require the internal program states to be identical during execution.

**Example.** Figure 3 shows an equivalent CUDA program pair. Both compute matrix-vector multiplication  $y = Ax$ , where  $A$  has dimensions  $(R, C)$  and  $x$  has size  $C$ . The right-hand program applies the *shared memory tiling* technique, loading  $x$  into shared memory  $tile$  (declared with `__shared__`). Synchronization primitives `__syncthreads()` are properly inserted to prevent synchronization issues.

**Automation.** The program transformation can be automated with tensor compilers, which provide a set of *schedules* to optimize loop-based programs. These schedules include loop tiling, loop fusion, loop reordering, loop unrolling, vectorization, and cache optimization. For any given schedule, the compiler can generate the transformed code. While different schedules can significantly impact program performance on the GPU, they do not affect the program’s correctness (assuming no compiler bugs), providing the foundation for automation.

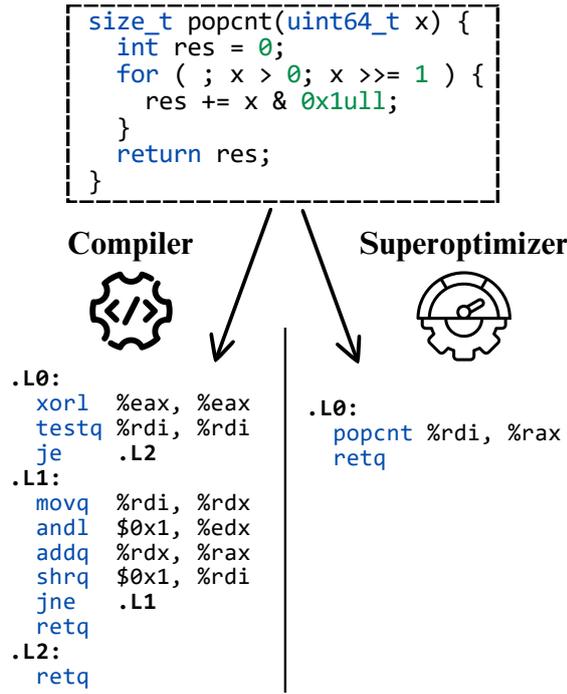


Figure 4: An equivalent pair from the x86-64 category in EquiBench. Both programs are compiled from the same C function shown above—the left using a compiler and the right using a *superoptimizer*. The function counts the number of set bits in the input `%rdi` register and stores the result in `%rax`. Their equivalence has been formally verified by the superoptimizer.

**Dataset Generation.** We utilize TVM as the tensor compiler (Chen et al., 2018) and sample tensor program schedules from TenSet (Zheng et al., 2021) to generate equivalent CUDA program pairs. Inequivalent pairs are created by sampling code from different tensor programs.

### 3.3 Pairs from a Superoptimizer (x86-64)

**Definition of Equivalence.** Two x86-64 assembly programs are considered equivalent if, for any input provided in the specified input registers, both programs produce identical outputs in the specified output registers. Differences in other registers or memory are ignored for equivalence checking.

**Example.** Figure 4 shows an example of an equivalent program pair in x86-64 assembly. Both programs implement the same C function, which counts the number of bits set to 1 in the variable `x` (mapped to the `%rdi` register) and stores the result in `%rax`. The left-hand program, generated by GCC with O3 optimization, uses a loop to count each bit individually, while the right-hand program, produced by a superoptimizer, leverages the `popcnt`

instruction, a hardware-supported operation for efficient bit counting. The superoptimizer verifies that both programs are semantically equivalent. Determining this equivalence requires a solid understanding of x86-64 assembly semantics and the ability to reason about all possible bit patterns.

**Automation.** A superoptimizer searches a space of programs to find one equivalent to the target. Test cases efficiently prune incorrect candidates, while formal verification guarantees the correctness of the optimized program. Superoptimizers apply aggressive and non-local transformations, making semantic equivalence reasoning more challenging. For example, in Figure 4, while a traditional compiler translates the loop in the source C program into a loop in assembly, a superoptimizer can find a more optimal instruction sequence by leveraging specialized hardware instructions. Such semantic equivalence is beyond the scope of traditional compilers.

**Dataset Generation.** We use Stoke (Schkufza et al., 2013) to generate program pairs. Assembly programs are sampled from prior work (Koenig et al., 2021), and Stoke applies transformations to produce candidate programs. If verification succeeds, the pair is labeled as equivalent; if the generated test cases fail, it is labeled as inequivalent.

### 3.4 Pairs from Programming Contests

**Definition of Equivalence.** Two programs are considered equivalent if they solve the same problem by producing the same output for any valid input, as defined by the problem description. Both programs, along with the problem description, are provided to determine equivalence.

**Example.** Given the problem description in Figure 5, all four programs are equivalent as they correctly compute the Fibonacci number. The **OJ\_A** pairs demonstrate **algorithmic** equivalence—the left-hand program uses recursion, while the right-hand program employs a for-loop. The **OJ\_V** pairs are generated through **variable renaming**, a **pure syntactic transformation** that can obscure the program’s semantics by removing meaningful variable names. The **OJ\_VA** pairs combine **both** algorithmic differences and variable renaming.

**Dataset Generation.** We sample Python submissions using a publicly available dataset from Online Judge (OJ) (Puri et al., 2021). For OJ\_A pairs, accepted submissions are treated as equivalent, while

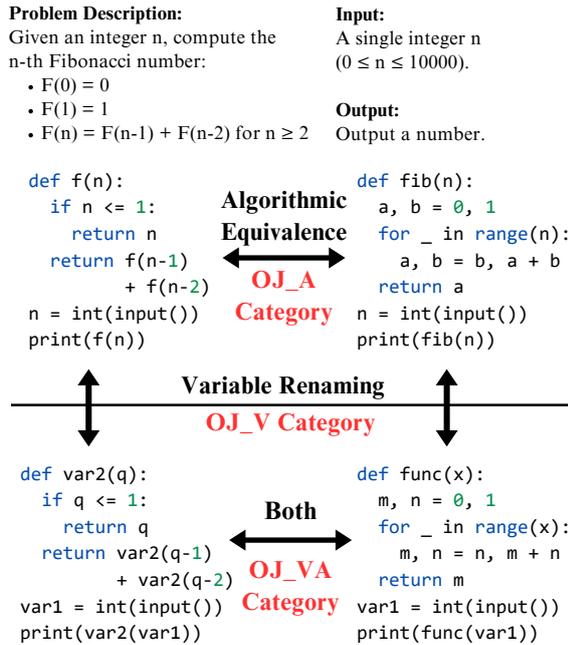


Figure 5: **Equivalent pairs from the OJ\_A, OJ\_V, OJ\_VA categories in EquiBench.** OJ\_A pairs demonstrate *algorithmic equivalence*, OJ\_V pairs involve *variable renaming* transformations, and OJ\_VA pairs combine *both* types of variations.

pairs consisting of an accepted submission and a wrong-answer submission are considered inequivalent. Variable renaming transformations are automated with an open-source tool (Flook, 2025).

## 4 Experimental Setup

**EquiBench.** Our dataset, EquiBench, consists of 2,400 program pairs across six equivalence categories. Each category contains 200 equivalent and 200 inequivalent pairs. Table 1 summarizes the lines of code, including the minimum, maximum, and average, for programs in each category, reflecting the wide variation in program lengths. As the dataset generation pipeline is fully automated, additional pairs can be generated as needed.

Category	Language	# Pairs	Lines of Code		
			Min	Max	Avg.
DCE	C	400	98	880	541
CUDA	CUDA	400	46	1733	437
x86-64	x86-64	400	8	29	14
OJ_A	Python	400	3	3403	82
OJ_V	Python	400	2	4087	70
OJ_VA	Python	400	3	744	35

Table 1: **Statistics of the EquiBench dataset.**

**Research Questions.** We investigate: 1) how different models perform on equivalence checking (Section 5.1); 2) whether prompting techniques, such as few-shot learning (Brown et al., 2020) and Chain-of-Thought (Wei et al., 2022), can enhance performance (Section 5.2); and 3) whether model predictions exhibit bias when judging program equivalence.

**Models.** We evaluate 17 large language models. For open-source models, including Mixtral (Jiang et al., 2024), Llama (Touvron et al., 2023), Qwen (Bai et al., 2023), DeepSeek (Liu et al., 2024a), we use Together AI, a model serving framework. For closed-source models (e.g., GPT-4 (Achiam et al., 2023), Claude-3.5 (Anthropic, 2024)), we access them via their official APIs, using the default temperature setting.

**Prompts.** The 0-shot evaluation is conducted using the prompt “You are here to judge if two programs are semantically equivalent. Here equivalence means {definition}. [Program 1]: {code1} [Program 2]: {code2} Please only output the answer of whether the two programs are equivalent or not. You should only output Yes or No.” The definition of equivalence and the corresponding program pairs are provided for each category. Additionally, for the categories of OJ\_A, OJ\_V and OJ\_VA, the prompt also includes the problem description. The full prompts used in our experiments for each equivalence category are in Appendix A.1.

**Error Handling.** Some models occasionally fail to follow the instruction to “output Yes or No”. To address this issue, we use GPT-4o to parse model outputs. In cases where no result can be extracted, we randomly assign “Yes” or “No” as the model’s output. These errors are very rare in advanced models but occur more frequently in smaller models.

## 5 Results

### 5.1 Model Accuracy

Table 2 shows the accuracy results for 17 state-of-the-art large language models on EquiBench under zero-shot prompting. Our findings are as follows:

**Reasoning models achieve the highest performance, demonstrating a clear advantage over non-reasoning models.** As shown in Table 2, reasoning models such as OpenAI o3-mini, DeepSeek R1, and o1-mini significantly outperform all others in our evaluation. This further underscores the

Model	DCE	CUDA	x86-64	OJ_A	OJ_V	OJ_VA	Overall Accuracy
<i>Random Baseline</i>	50.0	50.0	50.0	50.0	50.0	50.0	50.0
Llama-3.2-3B-Instruct-Turbo	50.0	49.8	50.0	51.5	51.5	51.5	50.7
Llama-3.1-8B-Instruct-Turbo	41.8	49.8	50.5	57.5	75.5	56.8	55.3
Mistral-7B-Instruct-v0.3	51.0	57.2	73.8	50.7	50.5	50.2	55.6
Mixtral-8x7B-Instruct-v0.1	50.2	47.0	64.2	59.0	61.5	55.0	56.1
Mixtral-8x22B-Instruct-v0.1	46.8	49.0	62.7	63.5	76.0	62.7	60.1
Llama-3.1-70B-Instruct-Turbo	47.5	50.0	58.5	66.2	72.0	67.5	60.3
QwQ-32B-Preview	48.2	50.5	62.7	65.2	71.2	64.2	60.3
Qwen2.5-7B-Instruct-Turbo	50.5	49.2	58.0	62.0	80.8	63.0	60.6
gpt-4o-mini-2024-07-18	46.8	50.2	56.8	64.5	91.2	64.0	62.2
Qwen2.5-72B-Instruct-Turbo	42.8	56.0	64.8	72.0	76.5	70.8	63.8
Llama-3.1-405B-Instruct-Turbo	40.0	49.0	75.0	72.2	74.5	72.8	63.9
DeepSeek-V3	41.0	50.7	69.2	73.0	83.5	72.5	65.0
gpt-4o-2024-11-20	43.2	49.5	65.2	71.0	87.0	73.8	65.0
claude3.5-sonnet-2024-10-22	38.5	<b>62.3</b>	70.0	71.2	78.0	73.5	65.6
o1-mini-2024-09-12	55.8	50.7	74.2	80.0	89.8	78.8	71.5
DeepSeek-R1	52.2	61.0	78.2	79.8	<b>91.5</b>	78.0	73.5
o3-mini-2025-01-31	<b>68.8</b>	59.0	<b>84.5</b>	<b>84.2</b>	88.2	<b>83.2</b>	<b>78.0</b>
Mean	47.9	52.4	65.8	67.3	76.4	67.0	62.8

Table 2: **Accuracy of 17 models on EquiBench under 0-shot prompting.** We report accuracy for each of the six equivalence categories along with the overall accuracy.

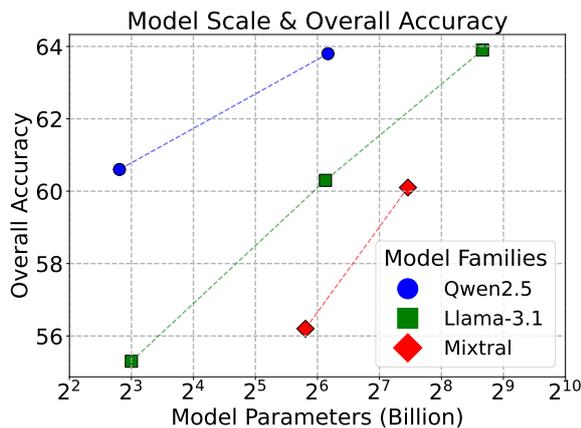


Figure 6: **Scaling Trend on EquiBench.**

complexity of equivalence checking as a code reasoning problem, where reasoning models exhibit a distinct advantage.

**EquiBench is a challenging benchmark.** Among the 17 models evaluated, OpenAI o3-mini achieves only 59.0% in the CUDA category despite being the top-performing model overall, with an accuracy of 78.0%. For the two most difficult categories, the highest accuracy across all models is 62.3% and 68.8%, respectively, only modestly above the random baseline of 50% accuracy for binary classification, highlighting the substantial room for improvement.

**Pure syntactic changes (OJ\_V) are the easiest for LLMs, while structural transformations are key to assessing deep semantic reasoning.** As

shown in the last row of Table 2, the OJ\_V category achieves the highest mean accuracy, with DeepSeek-R1 leading at 91.5%. This is because OJ\_V pairs are generated through trivial variable renaming, as seen in prior work (Badihi et al., 2021; Maveli et al., 2024). Additionally, combining variable renaming with algorithmic equivalence has little impact on difficulty, as indicated by the small drop in mean accuracy from OJ\_A 67.3% to OJ\_VA 67.0%. In contrast, all other categories involve non-local structural transformations, making them more challenging and essential for evaluating LLMs’ deep semantic reasoning.

### Scaling up models improves performance.

Larger models generally achieve better performance. Figure 6 shows scaling trends for the Qwen2.5, Llama-3.1, and Mixtral families, where accuracy improves with model size. The x-axis is on a logarithmic scale, highlighting how models exhibit consistent gains as parameters increase.

## 5.2 Prompting Strategies Analysis

We study few-shot in-context learning and Chain-of-Thought (CoT) prompting, evaluating four strategies: 0-shot, 4-shot, 0-shot with CoT, and 4-shot with CoT. For 4-shot, prompts include 2 equivalent and 2 inequivalent pairs. Appendix A.1 details the prompts, and Table 3 shows the results.

Our key finding is that **prompting strategies barely improve performance on EquiBench**, highlighting the task’s difficulty and need for

Model	0S	4S	0S-CoT	4S-CoT
o1-mini	71.5	71.5	<b>71.9</b>	<b>71.9</b>
gpt-4o	65.0	<b>66.5</b>	62.5	62.7
DeepSeek-V3	65.0	<b>66.9</b>	63.3	62.5
gpt-4o-mini	62.2	<b>63.5</b>	60.2	61.2

Table 3: **Accuracies of different prompting techniques.** We evaluate 0-shot and 4-shot in-context learning, both without and with Chain-of-Thought (CoT). Prompting strategies barely improve performance, highlighting the task’s difficulty and the need for task-specific approaches.

deeper reasoning. Few-shot prompting provides only minor improvements over 0-shot, while Chain-of-Thought shows slight benefits for o1-mini but marginally reduces performance for other models, underscoring the task’s complexity and the need for more advanced, task-specific approaches.

### 5.3 Bias in Model Prediction

We evaluate the prediction bias of the models and observe **a pronounced tendency to misclassify equivalent programs as inequivalent in the CUDA and x86-64 categories.** Table 4 presents the results for four representative models, showing high accuracy for inequivalent pairs but significantly lower accuracy for equivalent pairs, with full results for all models in Appendix A.2.

The bias in the CUDA category arises from extensive structural transformations, such as loop restructuring and shared memory optimizations, which make paired programs appear substantially different. In the x86-64 category, superoptimization applies non-local transformations to achieve optimal instruction sequences, introducing aggressive code restructuring that complicates equivalence reasoning and leads models to frequently misclassify equivalent pairs as inequivalent.

### 5.4 Case Studies

**Models lack capabilities for sound equivalence checking.** We find that simple changes that lead to semantic differences can confuse the models, causing them to produce incorrect predictions despite their correct predictions on the original program pairs. For example, o3-mini, which is one of the top-performing models in CUDA category, can correctly classify the pair shown in Figure 3 as equivalent. Next, we introduce synchronization bugs into the right-hand program, creating two inequivalent pairs with the original left-hand program: (1) removing the first `__syncthreads()`;

Model	CUDA		x86-64	
	Eq	Ineq	Eq	Ineq
<i>Random Baseline</i>	50.0	50.0	50.0	50.0
o3-mini	27.5	90.5	69.5	99.5
o1-mini	2.5	99.0	50.0	98.5
DeepSeek-R1	28.0	94.0	57.5	99.0
DeepSeek-V3	8.5	93.0	44.0	94.5

Table 4: Accuracies on equivalent and inequivalent pairs in the CUDA and x86-64 categories under 0-shot prompting, showing that **models perform significantly better on inequivalent pairs.** Random guessing serves as an unbiased baseline for comparison. Full results for all models are shown in Appendix A.2.

allows reads before all writes complete, causing race conditions; (2) removing the second `__syncthreads()`; lets faster threads overwrite shared data while slower threads read it. Despite these semantic differences, o3-mini misclassifies both pairs as equivalent.

**Proper hints enable models to correct misjudgments.** After o3-mini misclassifies the modified pairs, a hint about removed synchronization primitives allows it to correctly identify both as inequivalent, with accurate explanations highlighting data races. This suggests that training models on dedicated program analysis datasets, beyond only raw source code, may be useful for improving their code reasoning capabilities.

## 6 Conclusion

This paper presents EquiBench, a dataset for evaluating the code reasoning capabilities of large language models via program equivalence checking. Spanning four programming languages and six equivalence categories, EquiBench challenges models with diverse (in)equivalent program pairs generated through automated transformations, including syntactic changes, structural modifications, and algorithmic equivalence. Our evaluation shows that the best-performing model, OpenAI o3-mini, achieves only 59.0% in the CUDA category and 78.0% overall, with the most challenging categories achieving the best accuracies of just 62.3% and 68.8%, only modestly above the 50% random baseline. Few-shot learning and Chain-of-Thought prompting yield minimal gains, and models exhibit bias toward classifying programs with significant transformations as inequivalent. EquiBench provides a critical benchmark for advancing LLM-based code reasoning.

556  
557  
558  
559  
560  
561  
562  
563  
  
564  
  
565  
566  
567  
568  
569  
  
570  
571  
  
572  
573  
574  
575  
576  
  
577  
578  
579  
580  
581  
582  
583  
  
584  
585  
586  
587  
588  
  
589  
590  
591  
592  
  
593  
594  
595  
596  
597  
598  
  
599  
600  
601  
602  
603  
604  
  
605  
606  
607  
608

## Limitations

We make every effort to ensure that all pairs are correctly labeled, but cannot guarantee complete accuracy due to potential bugs in the toolchains or errors in the inputs (e.g., solutions from programming contests may be accepted based on a limited set of test cases that might not fully expose underlying bugs in the accepted solutions).

## References

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*.

Anthropic. 2024. Anthropic. <https://www.anthropic.com/news/claude-3-family>.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Sahar Badihi, Faridah Akinotcho, Yi Li, and Julia Rubin. 2020. Ardif: scaling program equivalence checking via iterative abstraction and refinement of common code. In *Proceedings of the 28th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*, pages 13–24.

Sahar Badihi, Yi Li, and Julia Rubin. 2021. Eqbench: A dataset of equivalent and non-equivalent program pairs. In *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*, pages 610–614. IEEE.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

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*, 33:1877–1901.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.

Junkai Chen, Zhiyuan Pan, Xing Hu, Zhenhao Li, Ge Li, and Xin Xia. 2024a. Reasoning runtime behavior of a program with llm: How far are we? *arXiv preprint cs.SE/2403.16437*.

Liangyu Chen, Bo Li, Sheng Shen, Jingkang Yang, Chunyuan Li, Kurt Keutzer, Trevor Darrell, and Zhiwei Liu. 2024b. Large language models are visual reasoning coordinators. *Advances in Neural Information Processing Systems*, 36.

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*.

Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. 2018. Tvm: An automated end-to-end optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 578–594.

Berkeley Churchill, Oded Padon, Rahul Sharma, and Alex Aiken. 2019. Semantic program alignment for equivalence checking. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 1027–1040.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Chris Cummins, Volker Seeker, Dejan Grubisic, Mostafa Elhoushi, Youwei Liang, Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Kim Hazelwood, Gabriel Synnaeve, et al. 2023. Large language models for compiler optimization. *arXiv preprint arXiv:2309.07062*.

Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Roziere, Jonas Gehring, Gabriel Synnaeve, and Hugh Leather. 2024. Meta large language model compiler: Foundation models of compiler optimization. *arXiv preprint arXiv:2407.02524*.

Manjeet Dahiya and Sorav Bansal. 2017. Black-box equivalence checking across compiler optimizations. In *Asian Symposium on Programming Languages and Systems*, pages 127–147. Springer.

Haoran Ding, Zhaoguo Wang, Yicun Yang, Dexin Zhang, Zhenglin Xu, Haibo Chen, Ruzica Piskac, and Jinyang Li. 2023. Proving query equivalence using linear integer arithmetic. *Proceedings of the ACM on Management of Data*, 1(4):1–26.

Yangruibo Ding, Jinjun Peng, Marcus J Min, Gail Kaiser, Junfeng Yang, and Baishakhi Ray. 2024.

665	Semcoder: Training code language models with comprehensive semantics reasoning. <i>arXiv preprint arXiv:2406.01006</i> .	718
666		719
667		720
668	Dennis Felsing, Sarah Grebing, Vladimir Klebanov, Philipp Rümmer, and Mattias Ulbrich. 2014. Automating regression verification. In <i>Proceedings of the 29th ACM/IEEE international conference on Automated software engineering</i> , pages 349–360.	721
669		722
670		723
671		724
672		725
673	Daniel Flook. 2025. Python variable renaming tool. <a href="https://github.com/dflook/python-minifier">https://github.com/dflook/python-minifier</a> .	726
674		727
675	Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. 2024. Cruxeval: A benchmark for code reasoning, understanding and execution. <i>arXiv preprint arXiv:2401.03065</i> .	728
676		729
677		730
678		731
679		732
680	Shubhani Gupta, Aseem Saxena, Anmol Mahajan, and Sorav Bansal. 2018. Effective use of smt solvers for program equivalence checking through invariant-sketching and query-decomposition. In <i>International Conference on Theory and Applications of Satisfiability Testing</i> , pages 365–382. Springer.	733
681		734
682		735
683		736
684		737
685		738
686	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022. Large language models are reasoning teachers. <i>arXiv preprint arXiv:2212.10071</i> .	739
687		740
688		741
689	Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. <i>arXiv preprint arXiv:2212.10403</i> .	742
690		743
691		744
692	Ali Reza Ibrahimzada, Kaiyao Ke, Mrigank Pawagi, Muhammad Salman Abid, Rangeet Pan, Saurabh Sinha, and Reyhaneh Jabbarvand. 2024. Repository-level compositional code translation and validation. <i>arXiv preprint arXiv:2410.24117</i> .	745
693		746
694		747
695		748
696		749
697	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Live-codebench: Holistic and contamination free evaluation of large language models for code. <i>arXiv preprint arXiv:2403.07974</i> .	750
698		751
699		752
700		753
701		754
702		755
703	Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	756
704		757
705		758
706		759
707		760
708	Jason R Koenig, Oded Padon, and Alex Aiken. 2021. Adaptive restarts for stochastic synthesis. In <i>Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation</i> , pages 696–709.	761
709		762
710		763
711		764
712		765
713	Emanuele La Malfa, Christoph Weinhuber, Orazio Torre, Fangru Lin, Samuele Marro, Anthony Cohn, Nigel Shadbolt, and Michael Wooldridge. 2024. Code simulation challenges for large language models. <i>arXiv preprint arXiv:2401.09074</i> .	766
714		767
715		768
716		769
717		770
		771
		772
		773
	Chris Lattner and Vikram Adve. 2004. Llvvm: A compilation framework for lifelong program analysis & transformation. In <i>International symposium on code generation and optimization, 2004. CGO 2004.</i> , pages 75–86. IEEE.	718
		719
		720
		721
		722
	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024a. Deepseek-v3 technical report. <i>arXiv preprint arXiv:2412.19437</i> .	723
		724
		725
		726
		727
	Changshu Liu, Shizhuo Dylan Zhang, Ali Reza Ibrahimzada, and Reyhaneh Jabbarvand. 2024b. Codemind: A framework to challenge large language models for code reasoning. <i>arXiv preprint arXiv:2402.09664</i> .	728
		729
		730
		731
		732
	Chenxiao Liu, Shuai Lu, Weizhu Chen, Daxin Jiang, Alexey Svyatkovskiy, Shengyu Fu, Neel Sundaresan, and Nan Duan. 2023. Code execution with pre-trained language models. <i>arXiv preprint arXiv:2305.05383</i> .	733
		734
		735
		736
		737
	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie LIU. 2021. CodeXGLUE: A machine learning benchmark dataset for code understanding and generation. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)</i> .	738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
	Nickil Maveli, Antonio Vergari, and Shay B Cohen. 2024. What can large language models capture about code functional equivalence? <i>arXiv preprint arXiv:2408.11081</i> .	749
		750
		751
		752
	Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. <i>arXiv preprint arXiv:2410.05229</i> .	753
		754
		755
		756
		757
	Federico Mora, Yi Li, Julia Rubin, and Marsha Chechik. 2018. Client-specific equivalence checking. In <i>Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering</i> , pages 441–451.	758
		759
		760
		761
		762
	Ansong Ni, Miltiadis Allamanis, Arman Cohan, Yinlin Deng, Kensen Shi, Charles Sutton, and Pengcheng Yin. 2024. Next: Teaching large language models to reason about code execution. <i>arXiv preprint arXiv:2404.14662</i> .	763
		764
		765
		766
		767
	Shankara Pailoor, Yuepeng Wang, and Işıl Dillig. 2024. Semantic code refactoring for abstract data types. <i>Proceedings of the ACM on Programming Languages</i> , 8(POPL):816–847.	768
		769
		770
		771
	Rangeet Pan, Ali Reza Ibrahimzada, Rahul Krishna, Divya Sankar, Lambert Pougues Wassi, Michele	772
		773

774	Merler, Boris Sobolev, Raju Pavuluri, Saurabh Sinha, and Reyhaneh Jabbarvand. 2024. Lost in translation: A study of bugs introduced by large language models while translating code. In <i>Proceedings of the IEEE/ACM 46th International Conference on Software Engineering</i> , pages 1–13.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	829
775			830
776			831
777			832
778			833
779			
780	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Sean Shi, Michael Choi, Anish Agrawal, Arnab Chopra, et al. 2025. Humanity’s last exam. <i>arXiv preprint arXiv:2501.14249</i> .	Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. 2023. Automated program repair in the era of large pre-trained language models. In <i>2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)</i> , pages 1482–1494. IEEE.	834
781			835
782			836
783			837
784	Ruchir Puri, David S Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, et al. 2021. Codenet: A large-scale ai for code dataset for learning a diversity of coding tasks. <i>arXiv preprint arXiv:2105.12655</i> .	Chenyuan Yang, Yinlin Deng, Runyu Lu, Jiayi Yao, Jiawei Liu, Reyhaneh Jabbarvand, and Lingming Zhang. 2024a. Whitefox: White-box compiler fuzzing empowered by large language models. <i>Proceedings of the ACM on Programming Languages</i> , 8(OOPSLA2):709–735.	839
785			840
786			841
787			842
788			843
789			844
790	Eric Schkufza, Rahul Sharma, and Alex Aiken. 2013. Stochastic superoptimization. <i>ACM SIGARCH Computer Architecture News</i> , 41(1):305–316.	Chenyuan Yang, Zijie Zhao, and Lingming Zhang. 2023. Kernelgpt: Enhanced kernel fuzzing via large language models. <i>arXiv preprint arXiv:2401.00563</i> .	845
791			846
792			847
793	Rahul Sharma, Eric Schkufza, Berkeley Churchill, and Alex Aiken. 2013. Data-driven equivalence checking. In <i>Proceedings of the 2013 ACM SIGPLAN international conference on Object oriented programming systems languages &amp; applications</i> , pages 391–406.	Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and understanding bugs in c compilers. In <i>Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation</i> , pages 283–294.	848
794			849
795			850
796			851
797			852
798	Alex Shypula, Pengcheng Yin, Jeremy Lacomis, Claire Le Goues, Edward Schwartz, and Graham Neubig. 2021. Learning to superoptimize real-world programs. <i>arXiv preprint arXiv:2109.13498</i> .	Zhen Yang, Fang Liu, Zhongxing Yu, Jacky Wai Keung, Jia Li, Shuo Liu, Yifan Hong, Xiaoxue Ma, Zhi Jin, and Ge Li. 2024b. Exploring and unleashing the power of large language models in automated code translation. <i>Proceedings of the ACM on Software Engineering</i> , 1(FSE):1585–1608.	853
799			854
800			855
801			856
802	Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob Gardner, Milad Hashemi, Graham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. 2023. Learning performance-improving code edits. <i>arXiv preprint arXiv:2302.07867</i> .	Dylan Zhang, Curt Tigges, Zory Zhang, Stella Biderman, Maxim Raginsky, and Talia Ringer. 2024. Transformer-based models are not yet perfect at learning to emulate structural recursion. <i>arXiv preprint arXiv:2401.12947</i> .	857
803			858
804			859
805			860
806			861
807			862
808	Rajat Singh and Srikanta Bedathur. 2024. Exploring the use of llms for sql equivalence checking. <i>arXiv preprint arXiv:2412.05561</i> .	Fuheng Zhao, Lawrence Lim, Ishtiyaque Ahmad, Divyakant Agrawal, and Amr El Abbadi. 2023. Llm-sql-solver: Can llms determine sql equivalence? <i>arXiv preprint arXiv:2312.10321</i> .	863
809			864
810			865
811	Theodoros Theodoridis, Manuel Rigger, and Zhendong Su. 2022. Finding missed optimizations through the lens of dead code elimination. In <i>Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems</i> , pages 697–709.		866
812			867
813			868
814			869
815			870
816			871
817	Zhao Tian, Honglin Shu, Dong Wang, Xuejie Cao, Yasutaka Kamei, and Junjie Chen. 2024. Large language models for equivalent mutant detection: How far are we? In <i>Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis</i> , pages 1733–1745.	Lianmin Zheng, Ruo Chen Liu, Junru Shao, Tianqi Chen, Joseph E Gonzalez, Ion Stoica, and Ameer Haj Ali. 2021. Tenset: A large-scale program performance dataset for learned tensor compilers. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)</i> .	872
818			873
819			874
820			875
821			876
822			877
823	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. <i>arXiv preprint arXiv:2302.13971</i> .	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. <i>arXiv preprint arXiv:2205.10625</i> .	878
824			879
825			
826			
827			
828			

## 880 A Appendix

### 881 A.1 Prompts

#### 882 A.1.1 DCE Category

883 We show the prompts for 0-shot, 4-shot, 0-shot CoT, 4-shot CoT settings.

884 **0-Shot.** You are here to judge if two C programs are semantically equivalent.

885 Here equivalence means that, when run on the same input, the two programs always have the same  
886 program state at all corresponding points reachable by program execution.

887 [Program 1]:

888 {program\_1\_code}

890 [Program 2]:

891 {program\_2\_code}

893 Please only output the answer of whether the two programs are equivalent or not. You should only  
894 output YES or NO.

897 **0-shot CoT.** You are here to judge if two C programs are semantically equivalent.

898 Here equivalence means that, when run on the same input, the two programs always have the same  
899 program state at all corresponding points reachable by program execution.

900 [Program 1]:

901 {program\_1\_code}

902 [Program 2]:

903 {program\_2\_code}

905 Please output the answer of whether the two programs are equivalent or not. You should output YES or  
906 NO in the end. Let's think step by step.

909 **4-shot.** You are here to judge if two C programs are semantically equivalent.

910 Here equivalence means that, when run on the same input, the two programs always have the same  
911 program state at all corresponding points reachable by program execution.

912 **[Example 1]:**

913 [Program 1]:

```
914  
915     int main() {  
916         int x = 0;  
917         if (false) {  
918             x = 1;  
919         }  
920         return 0;  
921     }
```

922 [Program 2]:

```
923  
924     int main() {  
925         int x = 0;  
926         if (true) {  
927             x = 1;  
928         }
```

return 0;	929
}	930
[Answer]: NO	931
	932
<b>[Example 2]:</b>	933
[Program 1]:	934
	935
int main() {	936
int x = 0;	937
if (false) {	938
x = 1;	939
}	940
return 0;	941
}	942
[Program 2]:	943
	944
int main() {	945
int x = 0;	946
return 0;	947
}	948
[Answer]: YES	949
	950
<b>[Example 3]:</b>	951
[Program 1]:	952
	953
char b[2];	954
static int c = 0;	955
int main() {	956
if (&b[0] == &c) {	957
c = 1;	958
}	959
return 0;	960
}	961
[Program 2]:	962
	963
char b[2];	964
static int c = 0;	965
int main() {	966
c = 1;	967
return 0;	968
}	969
[Answer]: NO	970
	971
<b>[Example 4]:</b>	972
[Program 1]:	973
	974
char b[2];	975
static int c = 0;	976
int main() {	977
if (&b[0] == &c) {	978
c = 1;	979
}	980
return 0;	981
}	982
[Program 2]:	983
	984

```
985     char b[2];
986     static int c = 0;
987     int main() {
988         return 0;
989     }
```

990 [Answer]: YES

991 [Program 1]:

```
992
993
994     {program_1_code}
```

995 [Program 2]:

```
996
997     {program_2_code}
```

998 Please only output the answer of whether the two programs are equivalent or not. You should only  
999 output YES or NO.

1000  
1001 **4-shot CoT.** You are here to judge if two C programs are semantically equivalent.

1002 Here equivalence means that, when run on the same input, the two programs always have the same  
1003 program state at all corresponding points reachable by program execution.

1004  
1005 **[Example 1]:**

1006 [Program 1]:

```
1007
1008     int main() {
1009         int x = 0;
1010         if (false) {
1011             x = 1;
1012         }
1013         return 0;
1014     }
```

1015 [Program 2]:

```
1016
1017     int main() {
1018         int x = 0;
1019         if (true) {
1020             x = 1;
1021         }
1022         return 0;
1023     }
```

1024 [Answer]: x = 1 in program 1 will not be executed, but x = 1 in program 2 will be executed, leading to  
1025 different program states.

1026 The answer is NO.

1027  
1028 **[Example 2]:**

1029 [Program 1]:

```
1030
1031     int main() {
1032         int x = 0;
1033         if (false) {
1034             x = 1;
1035         }
1036         return 0;
1037     }
```

[Program 2]: 1038  
 1039  

```

int main() {
    int x = 0;
    return 0;
}

```

 1040  
 1041  
 1042  
 1043

[Answer]: x = 1 in program 1 will not be executed, and this statement does not exist in program 2. 1044  
 Program states are always the same. 1045  
 The answer is YES. 1046

**[Example 3]:** 1047

[Program 1]: 1048  
 1049  
 1050  

```

char b[2];
static int c = 0;
int main() {
    if (&b[0] == &c) {
        c = 1;
    }
    return 0;
}

```

 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058

[Program 2]: 1059

```

char b[2];
static int c = 0;
int main() {
    c = 1;
    return 0;
}

```

 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066

[Answer]: The if statement in program 1 checks whether the memory address of b[0] equals c's address. 1067  
 c = 1 will not be executed in program 1, leading to a program state different from program 2. 1068  
 The answer is NO. 1069

**[Example 4]:** 1070

[Program 1]: 1071

```

char b[2];
static int c = 0;
int main() {
    if (&b[0] == &c) {
        c = 1;
    }
    return 0;
}

```

 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079  
 1080  
 1081

[Program 2]: 1082

```

char b[2];
static int c = 0;
int main() {
    return 0;
}

```

 1083  
 1084  
 1085  
 1086  
 1087  
 1088

[Answer]: The if statement in program 1 checks whether the memory address of b[0] equals c's address. 1089  
 c = 1 will not be executed in program 1, so the two programs always have the same states. 1090  
 The answer is YES. 1091

[Program 1]: 1092

1093

1094 {program\_1\_code}

1095 [Program 2]:

1096

1097 {program\_2\_code}

1098 Please output the answer of whether the two programs are equivalent or not. You should output YES or  
1099 NO in the end. Let's think step by step.

### 1100 A.1.2 CUDA Category

1101 We show the prompts for 0-shot and 4-shot CoT settings.

1102 **0-Shot.** You are here to judge if two CUDA programs are semantically equivalent.

1103 Here equivalence means that, when run on the same valid input, the two programs always compute the  
1104 same mathematical output (neglecting floating point rounding errors).

1105 [Program 1]:

1106 {program\_1\_code}

1107 [Program 2]:

1108 {program\_2\_code}

1109 Please only output the answer of whether the two programs are equivalent or not. You should only  
1110 output YES or NO.

1111

1112 **4-shot CoT.** You are here to judge if two CUDA programs are semantically equivalent.

1113 Here equivalence means that, when run on the same valid input, the two programs always compute the  
1114 same mathematical output (neglecting floating point rounding errors).

1115

1116 **[Example 1]:**

1117 [Program 1]:

1118

```
1119 __global__ void sgemv_naive(int M, int N, int K, float alpha,  
1120 const float *A, const float *B, float beta, float *C) {  
1121     const uint x = blockIdx.x * blockDim.x + threadIdx.x;  
1122     const uint y = blockIdx.y * blockDim.y + threadIdx.y;  
1123  
1124     if (x < M && y < N) {  
1125         float tmp = 0.0;  
1126         for (int i = 0; i < K; ++i) {  
1127             tmp += A[x * K + i] * B[i * N + y];  
1128         }  
1129         C[x * N + y] = alpha * tmp + beta * C[x * N + y];  
1130     }  
1131 }
```

1132 [Program 2]:

1133

```
1134 __global__ void sgemv_naive(int M, int N, int K, float alpha,  
1135 const float *A, const float *B, float beta, float *C) {  
1136     const uint x = blockIdx.x * blockDim.x + threadIdx.x;  
1137     const uint y = blockIdx.y * blockDim.y + threadIdx.y;  
1138  
1139     if (x < M && y < N) {  
1140         float tmp = 0.0;  
1141         for (int i = 0; i < K; ++i) {  
1142             tmp += A[x * K + i] * B[i * N + y];  
1143         }  
1144         C[x * N + y] = beta * tmp + alpha * C[x * N + y];  
1145     }  
1146 }
```

[Answer]: Program 1 computes  $C = \alpha*(A@B) + \beta*C$ , while Program 2 computes  $C = \beta*(A@B) + \alpha*C$ .

The answer is NO.

**[Example 2]:**

[Program 1]:

```
__global__ void sgemm_naive(int M, int N, int K, float alpha,
    const float *A, const float *B, float beta, float *C) {
    const uint x = blockIdx.x * blockDim.x + threadIdx.x;
    const uint y = blockIdx.y * blockDim.y + threadIdx.y;

    if (x < M && y < N) {
        float tmp = 0.0;
        for (int i = 0; i < K; ++i) {
            tmp += A[x * K + i] * B[i * N + y];
        }
        C[x * N + y] = alpha * tmp + beta * C[x * N + y];
    }
}
```

[Program 2]:

```
template <const uint BLOCKSIZE>
__global__ void sgemm_global_mem_coalesce(int M, int N,
    int K, float alpha, const float *A, const float *B,
    float beta, float *C) {
    const int cRow = blockIdx.x * BLOCKSIZE
        + (threadIdx.x / BLOCKSIZE);
    const int cCol = blockIdx.y * BLOCKSIZE
        + (threadIdx.x % BLOCKSIZE);

    if (cRow < M && cCol < N) {
        float tmp = 0.0;
        for (int i = 0; i < K; ++i) {
            tmp += A[cRow * K + i] * B[i * N + cCol];
        }
        C[cRow * N + cCol] = alpha * tmp
            + beta * C[cRow * N + cCol];
    }
}
```

[Answer]: Both programs compute  $C = \alpha*(A@B) + \beta*C$ .

Program 2 improves performance with global memory coalescing, which does not change computation results.

The answer is YES.

**[Example 3]:**

[Program 1]:

```
__global__ void sgemm_naive(int M, int N, int K, float alpha,
    const float *A, const float *B, float beta, float *C) {
    const uint x = blockIdx.x * blockDim.x + threadIdx.x;
    const uint y = blockIdx.y * blockDim.y + threadIdx.y;

    if (x < M && y < N) {
        float tmp = 0.0;
        for (int i = 0; i < K; ++i) {
            tmp += A[x * K + i] * B[i * N + y];
        }
        C[x * N + y] = alpha * tmp + beta * C[x * N + y];
    }
}
```

1208 [Program 2]:

1209

```
1210     template <const int BLOCKSIZE>
1211     __global__ void sgemm_shared_mem_block(int M, int N, int K,
1212     float alpha, const float *A, const float *B, float beta,
1213     float *C) {
1214         const uint cRow = blockIdx.x;
1215         const uint cCol = blockIdx.y;
1216
1217         __shared__ float As[BLOCKSIZE * BLOCKSIZE];
1218         __shared__ float Bs[BLOCKSIZE * BLOCKSIZE];
1219
1220         const uint threadCol = threadIdx.x % BLOCKSIZE;
1221         const uint threadRow = threadIdx.x / BLOCKSIZE;
1222
1223         A += cRow * BLOCKSIZE * K;
1224         B += cCol * BLOCKSIZE;
1225         C += cRow * BLOCKSIZE * N + cCol * BLOCKSIZE;
1226
1227         float tmp = 0.0;
1228         for (int bkIdx = 0; bkIdx < K; bkIdx += BLOCKSIZE) {
1229             As[threadRow * BLOCKSIZE + threadCol] =
1230                 A[threadRow * K + threadCol];
1231             Bs[threadRow * BLOCKSIZE + threadCol] =
1232                 B[threadRow * N + threadCol];
1233
1234             A += BLOCKSIZE;
1235             B += BLOCKSIZE * N;
1236
1237             for (int dotIdx = 0; dotIdx < BLOCKSIZE; ++dotIdx) {
1238                 tmp += As[threadRow * BLOCKSIZE + dotIdx] *
1239                     Bs[dotIdx * BLOCKSIZE + threadCol];
1240             }
1241         }
1242         C[threadRow * N + threadCol] = alpha * tmp
1243             + beta * C[threadRow * N + threadCol];
1244     }
```

1245 [Answer]: Both programs aim to compute  $C = \alpha \cdot (A @ B) + \beta \cdot C$ , but there are two synchronization  
1246 bugs in Program 2.

1247 Before entering the inner loop to compute tmp, there is no guarantee that the cache (As, Bs) is fully  
1248 populated by all threads.

1249 At the end of each iteration of bkIdx, faster threads may fetch the next block into the cache before slower  
1250 threads are done.

1251 The answer is NO.

1252

1253 [Example 4]:

1254 [Program 1]:

1255

```
1256     __global__ void sgemm_naive(int M, int N, int K, float alpha,
1257     const float *A, const float *B, float beta, float *C) {
1258         const uint x = blockIdx.x * blockDim.x + threadIdx.x;
1259         const uint y = blockIdx.y * blockDim.y + threadIdx.y;
1260
1261         if (x < M && y < N) {
1262             float tmp = 0.0;
1263             for (int i = 0; i < K; ++i) {
1264                 tmp += A[x * K + i] * B[i * N + y];
1265             }
1266             C[x * N + y] = alpha * tmp + beta * C[x * N + y];
1267         }
1268     }
```

1269 [Program 2]:

1270

```

template <const int BLOCKSIZE>
__global__ void sgemm_shared_mem_block(int M, int N, int K,
float alpha, const float *A, const float *B, float beta,
float *C) {
const uint cRow = blockIdx.x;
const uint cCol = blockIdx.y;

__shared__ float As[BLOCKSIZE * BLOCKSIZE];
__shared__ float Bs[BLOCKSIZE * BLOCKSIZE];

const uint threadCol = threadIdx.x % BLOCKSIZE;
const uint threadRow = threadIdx.x / BLOCKSIZE;

A += cRow * BLOCKSIZE * K;
B += cCol * BLOCKSIZE;
C += cRow * BLOCKSIZE * N + cCol * BLOCKSIZE;

float tmp = 0.0;
for (int bkIdx = 0; bkIdx < K; bkIdx += BLOCKSIZE) {
As[threadRow * BLOCKSIZE + threadCol] =
A[threadRow * K + threadCol];
Bs[threadRow * BLOCKSIZE + threadCol] =
B[threadRow * N + threadCol];

__syncthreads();
A += BLOCKSIZE;
B += BLOCKSIZE * N;

for (int dotIdx = 0; dotIdx < BLOCKSIZE; ++dotIdx) {
tmp += As[threadRow * BLOCKSIZE + dotIdx] *
Bs[dotIdx * BLOCKSIZE + threadCol];
}
__syncthreads();
}
C[threadRow * N + threadCol] = alpha * tmp
+ beta * C[threadRow * N + threadCol];
}

```

[Answer]: Both programs aim to compute  $C = \alpha(A@B) + \beta * C$ . 1308

Program 2 load a chunk of A and a chunk of B from global memory into shared memory. 1309

Such shared memory cache-blocking improves performance but does not change the correctness of the computation (no bugs found). 1310

The answer is YES. 1312

[Program 1]: 1313

```
{program_1_code} 1314
```

[Program 2]: 1316

```
{program_2_code} 1317
```

Please output the answer of whether the two programs are equivalent or not. You should output YES or NO in the end. Let's think step by step. 1319

### A.1.3 x86-64 Category 1321

We show the prompts for 0-shot and 4-shot CoT settings. 1322

**0-shot.** You are here to judge if two x86-64 programs are semantically equivalent. 1323

Here equivalence means that, given any input bits in the register {def\_in}, the two programs always have 1324

the same bits in register {live\_out}. Differences in other registers do not matter for equivalence checking. 1325

[Program 1]: 1326

1327

1328

1329 {program\_1\_code}

1330 [Program 2]:

1331

1332 {program\_2\_code}

1333 Please only output the answer of whether the two programs are equivalent or not. You should only  
1334 output YES or NO.

1335

1336 **4-shot CoT.** You are here to judge if two x86-64 programs are semantically equivalent.

1337 Here equivalence means that, given any input bits in the register {def\_in}, the two programs always have  
1338 the same bits in register {live\_out}. Differences in other registers do not matter for equivalence checking.

1339

1340 **[Example 1]:** In this example, the input register is %rdi, and output register is %rdi.

1341 [Program 1]:

1342

```
1343 movq -8(%rsp), %rdi
1344 .L4:
1345 call (%rdi)
1346 movq 8(%rdi), %rdi
1347 .L6:
1348 testq %rdi, %rdi
1349 jne .L4
```

1350 [Program 2]:

1351

```
1352 .L4:
1353 movq -8(%rsp), %rdi
1354 call (%rdi)
1355 movq 8(%rdi), %rdi
1356 movq %rdi, -8(%rsp)
1357 .L6:
1358 movq -8(%rsp), %rdi
1359 testq %rdi, %rdi
1360 jne .L4
```

1361 [Answer]: The additional instructions in Program 2 are: `movq %rdi, -8(%rsp)` and `movq -8(%rsp),`  
1362 `%rdi`.

1363 Program 2 stores the updated %rdi value back into -8(%rsp) after each iteration and reloads it before the  
1364 next iteration. But this does not affect the value of %rdi.

1365 The answer is YES.

1366

1367 **[Example 2]:** In this example, the input register is %rdi, and output register is %rdi.

1368 [Program 1]:

1369

```
1370 movq -8(%rsp), %rdi
1371 .L4:
1372 call (%rdi)
1373 movq 8(%rdi), %rdi
1374 .L6:
1375 testq %rdi, %rdi
1376 jne .L4
```

1377 [Program 2]:

1378

```
1379 .L4:
1380 movq -8(%rsp), %rdi
1381 call (%rdi)
1382 movq 8(%rdi), %rdi
1383 movq %rdi, -8(%rsp)
```

```

.L6:
movq -8(%rsp), %rdi
addq $1, %rdi
testq %rdi, %rdi
jne .L4

```

[Answer]: The additional instruction from Program 2 includes `addq $1, %rdi`, which increments `%rdi` by 1 before the test condition.

The two programs do not produce the same result for `%rdi`.

The answer is NO.

**[Example 3]:** In this example, the input register is `%rdi`, and output register is `%rax`.

[Program 1]:

```

.text
.globl _Z6popcntm
.type _Z6popcntm, @function
_Z6popcntm:
xorl %eax,%eax
testq %rdi,%rdi
je .L_4005b0
nop
.L_4005a0:
movq %rdi,%rdx
andl $0x1,%edx
addq %rdx,%rax
shrq $0x1,%rdi
jne .L_4005a0
retq
.L_4005b0:
retq
nop
nop
.size _Z6popcntm, .-_Z6popcntm

```

[Program 2]:

```

.text
.globl _Z6popcntm
.type _Z6popcntm @function
_Z6popcntm:
popcnt %rdi, %rax
retq
.size _Z6popcntm, .-_Z6popcntm

```

[Answer]: Both programs compute the population count (the number of 1s in a number's binary representation) of `%rdi` and store the result in `%rax`.

The answer is YES.

**[Example 4]:** In this example, the input register is `%rdi`, and output register is `%rax`.

[Program 1]:

```

.text
.globl _Z6popcntm
.type _Z6popcntm, @function
_Z6popcntm:
xorl %eax, %eax
testq %rdi, %rdi
je .L_4005b0
nop
.L_4005a0:
movq %rdi, %rdx
andl $0x1, %edx
addq %rdx, %rax

```

```

1445     addq    $1, %rax
1446     shrq    $0x1, %rdi
1447     jne    .L_4005a0
1448     retq
1449     .L_4005b0:
1450     retq
1451     nop
1452     nop
1453     .size  _Z6popcntm, .-_Z6popcntm

```

1454 [Program 2]:

```

1455
1456     .text
1457     .globl _Z6popcntm
1458     .type  _Z6popcntm @function
1459     _Z6popcntm:
1460     popcnt %rdi, %rax
1461     retq
1462     .size  _Z6popcntm, .-_Z6popcntm

```

1463 [Answer]: The instruction `addq $1, %rax` in Program 1 introduces a discrepancy by adding the number  
1464 of loop iterations to the output register.

1465 Program 2 simply computes the population count, but Program 1 adds an extra increment for each bit in  
1466 `%rdi`.

1467 The answer is NO.

1468  
1469 The input register is `{def_in}`, and the output register is `{live_out}`.

1470 [Program 1]:

```

1471
1472 {program_1_code}

```

1473 [Program 2]:

```

1474
1475 {program_2_code}

```

1476 Please output the answer of whether the two programs are equivalent or not. You should output YES or  
1477 NO in the end. Let's think step by step.

#### 1478 **A.1.4 OJ\_A Category**

1479 We show the prompts for both 0-shot and 4-shot CoT settings.

1480 **0-shot.** You are here to judge if two Python programs are semantically equivalent.

1481 You will be given [Problem Description], [Program 1] and [Program 2].

1482 Here equivalence means that, given any valid input under the problem description, the two programs will  
1483 always give the same output.

1484  
1485 [Problem Description]:

```

1486
1487 {problem_html}

```

1488 [Program 1]:

```

1489
1490 {program_1_code}

```

1491 [Program 2]:

```

1492
1493 {program_2_code}

```

1494 Please only output the answer of whether the two programs are equivalent or not. You should only  
1495 output YES or NO.

**4-shot CoT.** You are here to judge if two Python programs are semantically equivalent. 1496  
You will be given [Problem Description], [Program 1], and [Program 2]. 1497  
Here equivalence means that, given any valid input under the problem description, the two programs will 1498  
always give the same output. 1499

**[Example 1]:** 1501

[Problem Description]: 1502

Given a single line of input containing integers separated by spaces, sort the integers in ascending order 1503  
and print them in a single line separated by spaces. 1504

Input: A single line containing integers  $A[i]$  ( $-10^6 \leq A[i] \leq 10^6$ ,  $1 \leq n \leq 10^6$ ). 1505

Output: A single line of integers sorted in ascending order. 1506

Example Input: 4 2 5 1 3 1507

Example Output: 1 2 3 4 5 1508

[Program 1]: 1509

```
def bubble_sort(arr): 1510
    n = len(arr) 1511
    for i in range(n - 1): 1512
        for j in range(n - 1 - i): 1513
            if arr[j] > arr[j + 1]: 1514
                arr[j], arr[j + 1] = arr[j + 1], arr[j] 1515
    return arr 1516

nums = list(map(int, input().split())) 1517
sorted_nums = bubble_sort(nums) 1518
print("_".join(map(str, sorted_nums))) 1519
```

[Program 2]: 1520

```
def insertion_sort(arr): 1521
    for i in range(1, len(arr)): 1522
        key = arr[i] 1523
        j = i - 1 1524
        while j >= 0 and arr[j] > key: 1525
            arr[j + 1] = arr[j] 1526
            j -= 1 1527
        arr[j + 1] = key 1528
    return arr 1529

nums = list(map(int, input().split())) 1530
sorted_nums = insertion_sort(nums) 1531
print("_".join(map(str, sorted_nums))) 1532
```

[Answer]: Program 1 is bubble sort, and Program 2 is insertion sort. 1533  
The answer is YES. 1534

**[Example 2]:** 1535

[Problem Description]: Same as Example 1. 1536

[Program 1]: Same as Program 1 from Example 1. 1537

[Program 2]: 1538

```
def insertion_sort(arr): 1539
    for i in range(1, len(arr)): 1540
        key = arr[i] 1541
        j = i - 1 1542
        while j >= 0 and arr[j] < key: 1543
            arr[j + 1] = arr[j] 1544
            j -= 1 1545
        arr[j + 1] = key 1546
```

```

1554         return arr
1555
1556     nums = list(map(int, input().split()))
1557     sorted_nums = insertion_sort(nums)
1558     print("_".join(map(str, sorted_nums)))

```

1559 [Answer]: Program 1 is bubble sort, and Program 2 has a bug (the loop condition incorrectly uses arr[j]  
1560 < key instead of arr[j] > key).

1561 The answer is NO.

1562

1563 **[Example 3]:**

1564 [Problem Description]: Same as Example 1.

1565 [Program 1]:

1566

```

1567 def bubble_sort(arr):
1568     n = len(arr)
1569     for i in range(n - 1):
1570         for j in range(n - 1 - i):
1571             if arr[j] < arr[j + 1]:
1572                 arr[j], arr[j + 1] = arr[j + 1], arr[j]
1573     return arr
1574
1575     nums = list(map(int, input().split()))
1576     sorted_nums = bubble_sort(nums)
1577     print("_".join(map(str, sorted_nums)))

```

1578 [Program 2]: Same as Program 2 from Example 1.

1579

1580 [Answer]: Program 1 has a bug for bubble sort (the comparison is reversed, causing incorrect swaps).  
1581 The answer is NO.

1582

1583 **[Example 4]:**

1584 [Problem Description]: Same as Example 1.

1585 [Program 1]: Same as Program 1 from Example 1.

1586 [Program 2]:

1587

```

1588     nums = list(map(int, input().split()))
1589     sorted_nums = sorted(nums)
1590     print("_".join(map(str, sorted_nums)))

```

1591 [Answer]: Program 1 is bubble sort, and Program 2 uses Python's built-in sorting implementation.  
1592 The answer is YES.

1593

1594 [Problem Description]:

1595

1596 {problem\_html}

1597 [Program 1]:

1598

1599 {program\_1\_code}

1600 [Program 2]:

1601

1602 {program\_2\_code}

1603 Please output the answer of whether the two programs are equivalent or not. You should output YES or  
1604 NO in the end. Let's think step by step.

### 1605 A.1.5 OJ\_V Category

1606 We show the prompt for 4-shot CoT settings.

**4-shot CoT.** You are here to judge if two Python programs are semantically equivalent. 1607  
You will be given [Problem Description], [Program 1] and [Program 2]. 1608  
Here equivalence means that, given any valid input under the problem description, the two programs will 1609  
always give the same output. 1610

**[Example 1]:** 1612

[Problem Description]: 1613

Given a single line of input containing integers separated by spaces, sort the integers in ascending order 1614  
and print them in a single line separated by spaces. 1615

Input: A single line containing integers  $A[i]$  ( $-10^6 \leq A[i] \leq 10^6, 1 \leq n \leq 10^6$ ). 1616

Output: A single line of integers sorted in ascending order. 1617

Example Input: 1618

4 2 5 1 3 1619

Example Output: 1620

1 2 3 4 5 1621

[Program 1]: 1622

```
nums = list(map(int, input().split())) 1624
sorted_nums = sorted(nums) 1625
print("_".join(map(str, sorted_nums))) 1626
```

[Program 2]: 1627

```
random_var1 = list(map(int, input().split())) 1629
random_var2 = sorted(random_var1) 1630
print("_".join(map(str, random_var2))) 1631
```

[Answer]: The only difference is in variable names, which do not affect the logic or output of the program. 1632

The answer is YES. 1633

**[Example 2]:** 1635

[Problem Description]: 1636

Same as Example 1. 1638

[Program 1]: 1639

```
nums = list(map(int, input().split())) 1641
sorted_nums = sorted(nums) 1642
print("_".join(map(str, sorted_nums))) 1643
```

[Program 2]: 1644

```
nums = list(map(int, input().split())) 1646
sorted_nums = nums.sort() 1647
print("_".join(map(str, sorted_nums))) 1648
```

[Answer]: Program 1 sorts the integers in the correct way. In Program 2, `nums.sort()` modifies the list in 1649  
place and returns None. Program 2 will trigger a `TypeError`. 1650

The answer is NO. 1651

**[Example 3]:** 1653

[Problem Description]: 1654

Given a list of integers, remove all duplicate values while maintaining the order of their first appearance 1655  
and print the resulting list in a single line, separated by spaces. 1656

Input: A single line containing integers  $A[i]$  ( $-10^6 \leq A[i] \leq 10^6, 1 \leq n \leq 10^5$ ). 1657

Output: A single line containing the integers from the input with duplicates removed, in the order of their first appearance.

Example Input:

4 5 4 2 5 1 3

Example Output:

4 5 2 1 3

[Program 1]:

```
nums = list(map(int, input().split()))
unique_nums = []
for num in nums:
    if num not in unique_nums:
        unique_nums.append(num)
print("_".join(map(str, unique_nums)))
```

[Program 2]:

```
random_var1 = list(map(int, input().split()))
random_var2 = []
for random_var3 in random_var1:
    if random_var3 not in random_var2:
        random_var2.append(random_var3)
print("_".join(map(str, random_var2)))
```

[Answer]: The only difference is in variable names, which do not affect the logic or output of the program. The answer is YES.

#### [Example 4]:

[Problem Description]:

Same as Example 3.

[Program 1]:

```
nums = list(map(int, input().split()))
unique_nums = []
for num in nums:
    if num not in unique_nums:
        unique_nums.append(num)
print("_".join(map(str, unique_nums)))
```

[Program 2]:

```
nums = list(map(int, input().split()))
unique_nums = []
for num in nums:
    if num in unique_nums:
        unique_nums.append(num)
print("_".join(map(str, unique_nums)))
```

[Answer]: Program 1 correctly appends unique values to unique\_nums by checking if num not in unique\_nums.

Program 2 is incorrect because it uses if num in unique\_nums, causing only duplicates to be appended to the list.

The answer is NO.

[Problem Description]:

{problem\_html}

[Program 1]:	1712
	1713
{program_1_code}	1714
[Program 2]:	1715
	1716
{program_2_code}	1717
Please output the answer of whether the two programs are equivalent or not. You should output YES or NO in the end. Let's think step by step.	1718
	1719
<b>A.1.6 OJ_VA Category</b>	1720
We show the prompt for 4-shot CoT settings.	1721
<b>4-shot CoT.</b> You are here to judge if two Python programs are semantically equivalent.	1722
You will be given [Problem Description], [Program 1] and [Program 2].	1723
Here equivalence means that, given any valid input under the problem description, the two programs will always give the same output.	1724
	1725
	1726
<b>[Example 1]:</b>	1727
[Problem Description]:	1728
Given a single line of input containing integers separated by spaces, sort the integers in ascending order and print them in a single line separated by spaces.	1729
Input: A single line containing integers $A[i]$ ( $-10^6 \leq A[i] \leq 10^6$ , $1 \leq n \leq 10^6$ ).	1731
Output: A single line of integers sorted in ascending order.	1732
Example Input:	1733
4 2 5 1 3	1734
Example Output:	1735
1 2 3 4 5	1736
[Program 1]:	1737
	1738
<pre>def bubble_sort(arr):</pre>	1739
<pre>    n = len(arr)</pre>	1740
<pre>    for i in range(n - 1):</pre>	1741
<pre>        for j in range(n - 1 - i):</pre>	1742
<pre>            if arr[j] &gt; arr[j + 1]:</pre>	1743
<pre>                arr[j], arr[j + 1] = arr[j + 1], arr[j]</pre>	1744
<pre>    return arr</pre>	1745
	1746
<pre>nums = list(map(int, input().split()))</pre>	1747
<pre>sorted_nums = bubble_sort(nums)</pre>	1748
<pre>print("_".join(map(str, sorted_nums)))</pre>	1749
[Program 2]:	1750
	1751
<pre>def random_sort(rand_var1):</pre>	1752
<pre>    for rand_var2 in range(1, len(rand_var1)):</pre>	1753
<pre>        rand_var3 = rand_var1[rand_var2]</pre>	1754
<pre>        rand_var4 = rand_var2 - 1</pre>	1755
<pre>        while rand_var4 &gt;= 0 and rand_var1[rand_var4] &gt; rand_var3:</pre>	1756
<pre>            rand_var1[rand_var4 + 1] = rand_var1[rand_var4]</pre>	1757
<pre>            rand_var4 -= 1</pre>	1758
<pre>        rand_var1[rand_var4 + 1] = rand_var3</pre>	1759
<pre>    return rand_var1</pre>	1760
	1761
<pre>rand_input = list(map(int, input().split()))</pre>	1762
<pre>rand_output = random_sort(rand_input)</pre>	1763
<pre>print("_".join(map(str, rand_output)))</pre>	1764

1765 [Answer]: Program 1 is bubble sort, and Program 2 is insertion sort (though the variable names are  
1766 randomized).

1767 The answer is YES.

1768  
1769 **[Example 2]:**

1770 [Problem Description]:

1771 Same as Example 1.

1772 [Program 1]:

1773 Same as Program 1 from Example 1.

1774 [Program 2]:

1775  
1776 

```
def insertion_sort(rand_var1):  
1777     for i in range(1, len(rand_var1)):  
1778         key = rand_var1[i]  
1779         j = i - 1  
1780         while j >= 0 and rand_var1[j] < key:  
1781             rand_var1[j + 1] = rand_var1[j]  
1782             j -= 1  
1783         rand_var1[j + 1] = key  
1784     return rand_var1
```

1785  
1786 

```
nums = list(map(int, input().split()))  
1787 sorted_nums = insertion_sort(nums)  
1788 print("_".join(map(str, sorted_nums)))
```

1789 [Answer]: Program 1 is bubble sort, and Program 2 has a bug (the loop condition incorrectly uses  
1790 `rand_var1[j] < key` instead of `rand_var1[j] > key`).

1791 The answer is NO.

1792  
1793 **[Example 3]:**

1794 [Problem Description]:

1795 Same as Example 1.

1796 [Program 1]:

1797  
1798 

```
def rand_alg(rand_var):  
1799     n = len(rand_var)  
1800     for i in range(n - 1):  
1801         for j in range(n - 1 - i):  
1802             if rand_var[j] < rand_var[j + 1]:  
1803                 rand_var[j], rand_var[j + 1] = rand_var[j + 1], rand_var[j]  
1804     return rand_var
```

1805  
1806 

```
nums = list(map(int, input().split()))  
1807 sorted_nums = rand_alg(nums)  
1808 print("_".join(map(str, sorted_nums)))
```

1809 [Program 2]:

1810 Same as Program 2 from Example 1.

1811 [Answer]: Program 1 has a bug for bubble sort (the comparison is reversed, causing incorrect swaps).

1812 The answer is NO.

1813  
1814 **[Example 4]:**

1815 [Problem Description]:

1816 Same as Example 1.

1817 [Program 1]:

1818 Same as Program 1 from Example 1.

1819 [Program 2]:

1820

```
nums = list(map(int, input().split()))
sorted_nums = sorted(nums)
print("_".join(map(str, sorted_nums)))
```

[Answer]: Program 1 is bubble sort, and Program 2 uses Python's built-in sorting implementation.  
The answer is YES.

[Problem Description]:

{problem\_html}

[Program 1]:

{program\_1\_code}

[Program 2]:

{program\_2\_code}

Please output the answer of whether the two programs are equivalent or not. You should output YES or NO in the end. Let's think step by step.

1838  
1839  
1840  
1841

## A.2 Model Prediction Bias

We evaluate the prediction bias of the models and observe a pronounced tendency to misclassify equivalent programs as inequivalent in the CUDA and x86-64 categories. The table here shows the full results on all models under 0-shot prompting.

Model	CUDA		x86-64	
	Eq	Ineq	Eq	Ineq
<i>Random Baseline</i>	50.0	50.0	50.0	50.0
deepseek-ai/DeepSeek-V3	8.5	93.0	44.0	94.5
deepseek-ai/DeepSeek-R1	28.0	94.0	57.5	99.0
meta-llama/Llama-3.1-405B-Instruct-Turbo	6.0	92.0	68.5	81.5
meta-llama/Llama-3.1-8B-Instruct-Turbo	2.0	97.5	1.0	100.0
meta-llama/Llama-3.1-70B-Instruct-Turbo	7.0	93.0	27.5	89.5
meta-llama/Llama-3.2-3B-Instruct-Turbo	0.0	99.5	0.0	100.0
anthropic/claude-3-5-sonnet-20241022	62.5	62.0	49.5	90.5
Qwen/Qwen2.5-7B-Instruct-Turbo	18.5	80.0	17.5	98.5
Qwen/Qwen2.5-72B-Instruct-Turbo	14.5	97.5	36.0	93.5
Qwen/QwQ-32B-Preview	35.0	66.0	39.0	86.5
mistralai/Mixtral-8x7B-Instruct-v0.1	18.0	76.0	50.5	78.0
mistralai/Mixtral-8x22B-Instruct-v0.1	10.5	87.5	32.5	93.0
mistralai/Mistral-7B-Instruct-v0.3	52.5	62.0	87.0	60.5
openai/gpt-4o-mini-2024-07-18	0.5	100.0	16.5	97.0
openai/gpt-4o-2024-11-20	0.0	99.0	68.5	62.0
openai/o3-mini-2025-01-31	27.5	90.5	69.5	99.5
openai/o1-mini-2024-09-12	2.5	99.0	50.0	98.5