

TRACY: Benchmarking Execution Efficiency of LLM-Based Code Translation

Zhihao Gong¹, Zeyu Sun², Dong Huang³, Qingyuan Liang¹, Jie M. Zhang⁴, Dan Hao¹

¹Peking University, Beijing, China

²Institute of Software, Chinese Academy of Sciences, Beijing, China

³National University of Singapore, Singapore

⁴King’s College London, London, United Kingdom

zhihaogong@stu.pku.edu.cn, zeyu.zys@gmail.com, dong.huang@nus.edu.sg,
liangqy@stu.pku.edu.cn, jie.zhang@kcl.ac.uk, haodan@pku.edu.cn

Abstract

Automatic code translation is a fundamental task in modern software development. While the advent of Large Language Models (LLMs) has significantly improved the correctness of code translation, the critical dimension of execution efficiency remains overlooked. To address this gap, we introduce **TRACY**, the first comprehensive benchmark designed to evaluate the execution efficiency of LLM-translated code. TRACY is constructed through an LLM-driven two-stage pipeline: an initial stage generates a suite of stress tests to amplify performance differences, followed by an efficiency-oriented task pruning stage that isolates the efficiency-distinguishing tasks. The resulting benchmark comprises 1,011 code translation tasks across C++, Java, and Python, each accompanied by an average of 22.1 verified reference translations and 10 computationally demanding tests. Our extensive evaluation of 26 representative LLMs reveals that even top-tier LLMs struggle to consistently produce efficient code translations. For instance, *Claude-4-think*, the leading model for correctness, ranks eighth overall when time efficiency is taken into account, surpassed by several smaller open-source models. We further pinpoint that algorithmic flaws and improper resource handling are the most detrimental, causing a median time slowdown of 5.6× and memory increase of 12.0×, respectively. Our work underscores the necessity of jointly optimizing for correctness and efficiency in future LLM-based code translation.

1 Introduction

Automatic code translation, which refers to the process of migrating source code from one programming language to another, is a pivotal process in modern software engineering. It has facilitated critical scenarios such as modernizing legacy systems, extending library interoperability, and unifying disparate codebases under a standardized technology stack (Pan et al. 2024; Ibrahimzada et al. 2025).

The advent of Large Language Models (LLMs) has revolutionized the field of code intelligence (Zheng et al. 2023b), demonstrating profound capabilities in understanding and generating multilingual code (Zheng et al. 2023a; Jiang et al. 2024). Their ability to learn intricate programming patterns and semantic relationships from vast code corpora has significantly impacted various code-related tasks (Nijkamp et al. 2022; Seed et al. 2025), including code generation (Chen et al. 2021), completion (Husein, Aburajouh, and

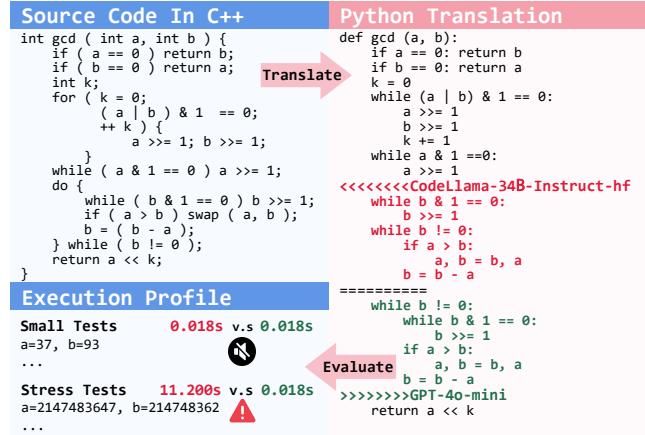


Figure 1: While functionally correct and indistinguishable on small tests, two LLM-generated code translations reveal an over 600× performance gap under computationally demanding stress tests due to a subtle algorithmic deviation.

Catal 2025), and debugging (Jimenez et al. 2023). Leveraging LLMs, the predominant research focus in code translation has been on enhancing the correctness of translated code. A surge of recent work explored techniques such as back translation (Ahmad et al. 2022), fine-tuning (Liu, Li, and Zhang 2023; Zhu et al. 2024), Chain-of-Thought (CoT) prompting (Macedo et al. 2024; Nitin, Krishna, and Ray 2024), and self-repair (Yang et al. 2024). These efforts have yielded code translations with demonstrably higher levels of syntactic and functional correctness.

Despite remarkable progress in correctness, the critical dimension of *code execution efficiency* (Niu et al. 2024; Vartziotis et al. 2024; Huang et al. 2024b) in LLM-translated code has been overlooked. In practice, efficiency directly impacts a software system’s responsiveness, scalability, and operational costs. A functionally correct but inefficient translation can introduce severe performance bottlenecks, rendering it unsuitable for real-world production deployment.

This issue is not a minor edge case and often evades detection by existing evaluation protocols. As illustrated in Figure 1, consider a C++ implementation of Stein’s GCD algorithm (Stein’s GCD 2025), a classic program-

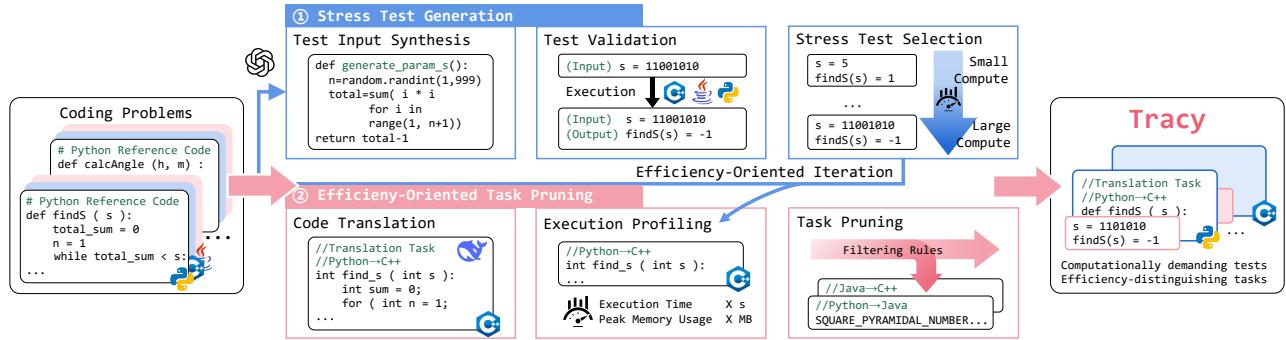


Figure 2: An overview of TRACY’s LLM-driven two-stage construction pipeline.

ming challenge from the widely-used *Transcoder-Test* code translation benchmark (Lachaux et al. 2020). We prompted two models (*GPT-4o-mini* and *CodeLlama-34B-Instructhf*) to translate this algorithm into Python. On the original tests from *Transcoder-Test*, both translations are performance indistinguishable. They successfully pass all the tests with negligible runtime differences (≈ 0.018 s). However, their performance diverges dramatically when challenged with more computationally demanding stress tests (e.g., $a=2,147,483,647$ and $b=2,147,483,62$). The translation from *CodeLlama* runs over $600\times$ slower than that from *GPT-4o-mini* (11.200s vs. 0.018s). Such a massive performance gap stems from a subtle structural misinterpretation, where the *CodeLlama*’s translation replaces efficient bitwise shifts with a slow and repetitive loop, thereby nullifying the algorithm’s core optimization of the source algorithm.

This example starkly demonstrates that functional equivalence does not guarantee efficiency equivalence in code translation, revealing a critical blind spot in current evaluation methodologies. Moreover, this flaw is pervasive, as our case study of 236 inefficient translations confirmed that the original benchmark’s tests overlooked 82.6% of performance degradation cases identified by more computationally demanding tests. The root of this blind spot is a singular focus on correctness, allowing functionally correct but performance-deficient code translation to pass the evaluation. Overcoming this limitation necessitates a new benchmark designed to assess efficiency alongside correctness.

To bridge the research gap, we introduce **TRACY** (**TRAnslated Code EfficiencY**), the first comprehensive benchmark for evaluating the execution efficiency of LLM-translated code. Drawing from the widely-used *Transcoder-Test*, TRACY meticulously curates 1,011 efficiency-critical code translation tasks from over 2,800 candidate tasks across C++, Java, and Python. The core of TRACY is an LLM-driven synthesis framework that generates a unified suite of computationally demanding stress tests. Unlike the original benchmark’s simple correctness checks, these stress tests are specifically designed to amplify latent performance differences. They increase the average execution time by over $9.4\times$ and peak memory usage by over $4.0\times$ even for the most efficient translations. This process effectively exposes algorithmic and idiomatic inefficiencies that simpler tests

would otherwise miss. Furthermore, to enable nuanced evaluation, each task is accompanied by an average of 22.1 verified reference translations, creating a robust spectrum for benchmarking the efficiency of candidate translations.

Leveraging TRACY, we conduct an extensive evaluation on 26 state-of-the-art LLMs. First, our end-to-end analysis uncovers a stark disparity between correctness and efficiency among LLM-translated code. For example, the correctness leader *Claude-4-think* passes 95.0% of tasks but ranks only eighth when considering time efficiency. A granular analysis further reveals that the disparity is highly dependent on the specific source and target languages involved in the translation. When translating from C++, models produce far more memory-efficient code when targeting Java than Python, with the average memory efficiency score being twice as high. Finally, our root-cause analysis on 236 inefficient cases pinpoints the inefficiency patterns: the most frequent inefficiencies result from the suboptimal idiom and library usage (61.9% of cases); in contrast, the most catastrophic degradation stems from algorithmic flaws and poor resource management, inducing a median slowdown of $5.6\times$ in time and an increase of $12.0\times$ in memory usage, respectively. This paper makes the following contributions:

1) Dimension. We spotlight efficiency as a critical yet overlooked dimension in LLM-based code translation. We introduce TRACY, the first benchmark for its evaluation.

2) Evaluation. We conduct an extensive evaluation of 26 state-of-the-art LLMs using TRACY, uncovering significant efficiency gaps and recurring inefficiency patterns.

3) Framework. We release the construction framework, the TRACY benchmark, and all evaluation results to foster further research.

2 The TRACY Benchmark

In the code translation context, a coding *problem* P represents a programming challenge (e.g., Stein’s GCD algorithm). A code translation task (abbr. *task*) is the process of translating a piece of *reference code* p from a source language l_s to the target language l_t ($p_{l_s} \rightarrow p_{l_t}$). We aim to construct a focused benchmark to highlight the crucial yet overlooked efficiency dimension of LLM-translated code.

We introduce TRACY, a new benchmark developed using an LLM-driven two-stage pipeline, as depicted in Fig-

ure 2. Our process begins with the high-quality problems sourced from Transcoder-Test (Lachaux et al. 2020; Yang et al. 2024). In the first stage, *stress test generation*, we generate a unified suite of challenging *stress tests* designed to magnify the potential performance disparities in code translations. In the second stage, *efficiency-oriented task pruning*, we systematically isolate code translation tasks where efficiency is a key distinguishing factor.

2.1 Problem Collection

TRACY sources its coding problems from the widely-used Transcoder-Test (Lachaux et al. 2020), which contains a diverse set of programming challenges originally from GeeksForGeeks. We opt to build upon this established foundation for two reasons: 1) It provides a high-quality and recognized set of problems, ensuring relevance and comparability with prior work; 2) More importantly, it allows us to directly highlight the overlooked efficiency dimension within the very same tasks previously studied only for correctness, encouraging future research to re-examine prior evaluation through the critical lens of efficiency. We build upon the most recent version of Transcoder-Test (Yang et al. 2024), which contains 568 problems and reference code in C++, Java, and Python, along with ten simple tests per problem.

2.2 Stress Test Generation

To create tests that can effectively amplify the potential efficiency gaps among LLM-based code translations, we employ an LLM-driven iterative stress test generation strategy, as detailed in Algorithm 1. This strategy progressively generates a unified suite of *stress tests*, which are defined as tests that demand substantial computational resources (e.g., execution time or peak memory usage) to exercise the code under test, making them highly effective at magnifying latent performance differences among code translations.

Test Input Synthesis. The initial step of stress test generation is to synthesize large test inputs sufficient to exercise the runtime behavior of functionally correct translations. As it is often infeasible to directly generate large-scale test inputs by an LLM due to its context window limitation, we utilize a *synthesis-based* approach (Lines 4-6) (Liu et al. 2024; Peng et al. 2025), where we prompt an advanced LLM (*GPT-4o*) to produce multiple test input *synthesizers* (\mathcal{S}) based on the reference code of the problem (available in C++, Java, and Python). Each synthesizer is a self-contained Python function that programmatically outputs a test input upon execution. To guide this process towards producing more computationally demanding test inputs, prompts in each iteration (Line 5) are enriched with few-shot examples (\mathcal{E}) of *synthesizers* that produce the most stressful tests in the previous iteration. This approach scalably bypasses the LLM’s context limit, enabling the generation of a broad range of large and randomized test inputs.

Test Validation. The test inputs from *synthesizers* undergo a rigorous validation process to ensure they are sound and consistent across all pieces of reference code (Lines 7-9). A test input is deemed valid for a problem only if it satisfies two criteria: (1) *Execution Integrity*, where it must execute successfully on all reference code (C++, Java, and

Algorithm 1: Stress Test Generation

```

1: FUNCTION GenerateStressTests( $P, I_{max}, K$ )
2:    $\mathcal{T}_s, \mathcal{E} \leftarrow \emptyset, \emptyset; iter \leftarrow 0$ 
3:   while  $iter < I_{max}$  do
4:     /* Stage 1: Test Input Synthesis */
5:      $\mathcal{P}_{prompt} \leftarrow \text{CreatePrompt}(P, \mathcal{E})$ 
6:      $\mathcal{S} \leftarrow \text{QueryLLM}(\mathcal{P}_{prompt})$ 
7:     /* Stage 2: Test Validation */
8:      $\mathcal{V} \leftarrow \bigcup_{s \in \mathcal{S}} \{\text{Execute}(s)\}$ 
9:      $\mathcal{T} \leftarrow \text{ValidateTest}(\mathcal{V}, P)$ 
10:    /* Stage 3: Stress Test Selection */
11:     $\mathcal{T}_{iter} \leftarrow \text{SelectTop}(\mathcal{T}, K, \text{time})$ 
12:     $\mathcal{T}_{iter} \leftarrow \mathcal{T}_{iter} \cup \text{SelectTop}(\mathcal{T} \setminus \mathcal{T}_{iter}, K, \text{mem})$ 
13:    /* Efficiency-Oriented Iteration */
14:    if HasConverged( $\mathcal{T}_{iter}, \mathcal{T}_s$ )
15:      then break
16:    end if
17:     $\mathcal{T}_s \leftarrow \text{UpdateTests}(\mathcal{T}_s, \mathcal{T}_{iter})$ 
18:     $\mathcal{E} \leftarrow \text{UpdateExamples}(\mathcal{T}_s)$ 
19:     $iter \leftarrow iter + 1$ 
20:  end while
21:  return  $\mathcal{T}_s$ 

```

Python) without raising runtime exceptions; and (2) *Output Consistency*, where it must produce identical outputs across all reference code. Our validation eliminates test inputs that are susceptible to language-specific behaviors (e.g., integer overflow in C++ versus Python’s arbitrary-precision integers), thereby establishing a reliable, language-agnostic ground truth for the resulting tests (\mathcal{T}).

Stress Test Selection. From the pool of validated tests, we select the most computationally demanding tests to form the set \mathcal{T}_{iter} for the current iteration (Lines 10-12). The selection process is based on two key efficiency indicators: *execution time* and *peak memory usage*, profiled for each test across all pieces of reference code. Specifically, we employ the *Borda Count* method (Emerson 2013) to create separate rankings for time and memory. For each indicator, a test’s score is the sum of its performance ranks (1st is the best) across C++, Java, and Python, where a lower score indicates a more resource-intensive test. Finally, the top- K tests from both the time and memory rankings are selected, collectively forming the most stressful tests for the current iteration.

Efficiency-Oriented Iteration. The entire process of stress test generation is designed as a feedback-driven loop to progressively discover more computationally demanding tests (Lines 13-18). This loop maintains a global pool (\mathcal{T}_s) of the most stressful tests found so far. After each iteration, the newly selected tests (\mathcal{T}_{iter}) are merged into this pool. The combined set is then re-ranked and pruned to keep only the overall top- K for both time and memory (Line 17). Synthesizers producing these top-performing tests are collected to serve as few-shot examples (\mathcal{E}) to enrich the prompt for the next iteration, steering the LLM towards generating even more computationally demanding test inputs (Line 18). The process terminates after I_{max} iterations or when it converges, defined as the point where no new tests are stressful enough to enter the global pool (Line 14).

Table 1: Comparison of TRACY with representative code translation benchmarks. *Effi.Dist.* means efficiency-distinguishable.

Benchmark	Source	#Probs.	Programming Languages	#Tests	Effi.Dist.	Evaluation
CodeNet (Puri et al. 2021)	Prog. Contests	200	C, C++, Go, Java, Python	1	✗	Correctness
Avatar (Ahmad et al. 2021)	Prog. Contests	250	Java, Python	25	✗	Correctness
TransCoder-Test (Roziere et al. 2021)	GeeksForGeeks	568	C++, Java, Python	10	✗	Correctness
G-TransEval (Jiao et al. 2023)	Multiple	400	C++, C#, Java, Python, JS	5	✗	Correctness
PolyHumanEval (Tao et al. 2024)	HumanEval	164	C++, C#, Java, Python (14 PLs)	7	✗	Correctness
TRACY (Ours)	GeeksForGeeks	361	C++, Java, Python	20 (10 stress)	✓	Efficiency Correctness

2.3 Efficiency-Oriented Task Pruning

Not all code translation tasks are suitable for the evaluation of efficiency. Intuitively, translating a trivial problem like *int add (int a, int b)* with LLMs is unlikely to yield variants with meaningful performance differences. We therefore apply a rigorous pruning process to ensure that every task retained in TRACY is genuinely efficiency-distinguishing.

For each task, we first generate a diverse pool of code translations, called *reference translations*, using a representative group of LLMs with sampling-based decoding. Next, all functionally correct reference translations are executed against our stress tests. We meticulously profile their performance by collecting key runtime indicators (execution time and peak memory usage). Finally, we apply three filtering criteria to each task based on these execution profiles: 1) *Feasibility*. We discard any task for which no evaluated LLM can produce a correct translation. These tasks offer no valid samples for efficiency profiling and are thus unsuitable. 2) *Impactfulness*. We discard simple tasks with trivial execution costs, where performance is not a practical concern. This is determined by requiring that at least one translation exceeds an execution time threshold (ϵ_T) or a peak memory usage threshold (ϵ_M) on stress tests. 3) *Diversity*. We discard tasks where correct translations exhibit negligible performance variance. This is quantified by requiring the coefficient of variation (CV) for either execution time or peak memory usage to be above a threshold (ϵ_D), ensuring the task can effectively differentiate models regarding the efficiency dimension. This stringent task filtering process ensures that every task included in TRACY presents a meaningful and distinguishable efficiency challenge.

2.4 Benchmark Implementation

The *stress test generation* stage utilizes *GPT-4o* with a temperature of 0.8 to generate synthesizers. For each problem, we prompt *GPT-4o* to generate 3 distinct test input synthesizers for each piece of reference code (C++, Java, and Python). Each of the resulting 9 synthesizers is then executed three times, producing a candidate pool of up to 27 unique inputs per iteration. For the iterative generation loop (Algorithm 1), we set a maximum of 5 iterations (I_{max} : 5) and selected the top $K = 5$ tests for both time and memory in each round. Consequently, the final stress test suite for each problem comprises 10 of the most computationally demanding tests identified across all the iterations.

For the *efficiency-oriented task pruning* stage, we first

generate a diverse pool of reference translations using 26 representative LLMs (See Section 3) with sampling-based decoding (temperature: 0.8). A translation is deemed correct only if it passes both the original Transcoder-Test’s tests and our new stress tests. Subsequently, we apply the filtering thresholds for impactfulness (ϵ_T : 0.001s, ϵ_M : 1.5MB) and diversity (ϵ_D : 0.05). Following prior work (Liu et al. 2024), the impactfulness thresholds (ϵ_T, ϵ_M) are empirically derived from the average resource consumption of a simple *Hello World* program across the three languages.

2.5 Benchmark Statistics

The final TRACY benchmark contains 361 coding problems, yielding a total of 1,011 efficiency-critical code translation tasks spanning six language pairs: C++ to Java (250), Python to Java (237), Java to Python (147), C++ to Python (140), Java to C++ (133), and Python to C++ (104).

As demonstrated in Table 1, TRACY introduces two fundamental departures from prior evaluation paradigms: 1) *Computationally Demanding Stress Tests*. TRACY is the first benchmark to integrate a dedicated suite of stress tests per problem, which are systematically generated to expose performance bottlenecks. These tests substantially magnify the average execution time of the most efficient LLM translations by 9.4× and their peak memory usage by 4.0×. The critical role of these tests is starkly highlighted by our findings: among 236 inefficient translations identified by our stress tests, the original Transcoder-Test’s simple tests fail to detect 82.6% inefficiency cases, flagging only 41. This disparity underscores the necessity of stress tests for a genuine assessment of efficiency. 2) *Efficiency-Distinguishing Tasks*. Unlike benchmarks that include all translation tasks from a source problem, TRACY exclusively contains tasks carefully selected through the rigorous pruning process, ensuring that every task in the benchmark presents a meaningful and distinguishable efficiency challenge. In addition, each task is equipped with a verified set of reference translations (avg. 22.1 per task), enabling nuanced evaluation against a known performance distribution instead of a single reference.

3 Evaluation Setup

We conduct an extensive evaluation using TRACY to assess the execution efficiency of LLM-translated code. We first outline models, metrics, and the evaluation protocol.

Models. We evaluate 26 representative Large Language Models (LLMs), spanning both leading proprietary LLMs

and prominent open-source code LLMs. For proprietary models, we include models from OpenAI (GPT-4o and O3 series), Anthropic (Claude-4-sonnet series), Google (Gemini-2.5 series), and DeepSeek. For open-source models, we assess the CodeLlama-hf (CL) (Roziere et al. 2023), DeepSeek-Coder (DSC) (Guo et al. 2024), and Qwen2.5-Coder (QC) series (Hui et al. 2024), including both base and instruction-tuned (Inst) variants. A more comprehensive model list is available in Appendix A.3.

Metrics. We assess model performance across two dimensions: *functional correctness* and *execution efficiency*.

For functional correctness, we report the *Passing Rate (Pass)* (Peng et al. 2025), which is the percentage of tasks for which the model’s translation passes all the tests.

For execution efficiency, we measure *execution time (Time)* and *peak memory usage (Mem.)*. To utilize a normalized score that takes efficiency into account, we adopt the *Beyond* score from Mercury (Du et al. 2024). This metric scores a candidate translation’s performance P (*Time* or *Mem.*) relative to the performance spectrum established by a set of verified translations (\mathcal{R}). For a correct solution, the score is computed as:

$$\text{Beyond}_P = \frac{\max(\mathcal{R}) - \text{clip}(P, \min(\mathcal{R}), \max(\mathcal{R}))}{\max(\mathcal{R}) - \min(\mathcal{R})} * 100\%$$

An incorrect translation automatically receives a Beyond score of 0%. In essence, a score of 100% indicates performance on par with or better than the most efficient in \mathcal{R} , while a score of 0% indicates performance equal to or worse than the least efficient. The final metrics are averaged across all the tasks for *Time* (B_T) and *Mem.* (B_M), respectively.

Evaluation Protocol. For each task, we prompt the model to output a single translation via *greedy decoding* (temperature=0.0) with zero-shot prompting. Its functional correctness is first validated against the original tests from Transcoder-Test. Any translation failing to compile, producing a runtime error, or exceeding the computational resource limits (180-second timeout, 4096-MB peak memory limit) is discarded as incorrect. The correct translation proceeds to the efficiency profiling stage, where it is executed on the stress tests to measure the execution time and peak memory usage. These figures are then scored for the Beyond metric. To mitigate measurement noise and ensure stability, we compute each indicator (*Time* and *Mem.*) as the arithmetic mean of five independent runs.

4 Evaluation

Our evaluation proceeds through three lenses: 1) An *End-to-end Evaluation* of overall model performance; 2) A *Granular Analysis* across different translation directions; 3) A *Root Cause Characterization* of common inefficiency patterns.

4.1 End-to-end Evaluation

We begin by quantifying the overall performance of the studied LLMs. Table 2 summarizes the results, revealing a significant disconnect between functional correctness and execution efficiency of LLM-translated code.

Table 2: The overall performance across LLMs. B_T/B_M : average Beyond scores (*Time/Mem.*) over all tasks. B_T^{com}/B_M^{com} : average Beyond score on the common subset of correct tasks. B_T^P/B_M^P : average Beyond score on each model’s *passed* translations only.

Model	Pass	B_T	B_T^{com}	B_T^P	B_M	B_M^{com}	B_M^P
Claude-4-think	95.0	49.8	49.9	52.4	48.7	48.9	51.3
Claude-4	94.3	47.5	48.0	50.4	46.7	46.1	49.5
DS-reasoner	72.1	39.0	-	54.1	24.6	-	34.1
DS-chat	89.4	43.2	46.1	48.3	42.9	45.6	48.0
Gemini-pro	62.0	37.3	-	60.2	22.0	-	35.4
Gemini-flash	61.6	36.8	-	59.7	22.6	-	36.4
GPT-4o	88.4	45.3	50.4	51.2	42.2	46.4	47.8
GPT-4o-mini	88.6	45.6	50.1	51.4	42.8	46.5	48.3
O3	84.1	49.4	-	58.8	31.3	-	37.3
O3-mini	89.9	42.4	45.9	47.2	49.5	53.1	55.0
CL-7B	61.3	35.8	-	58.4	22.7	-	37.0
CL-7B-Inst	74.8	44.2	-	59.1	27.7	-	37.0
CL-13B	60.0	35.9	-	59.8	21.9	-	36.4
CL-13B-Inst	76.2	44.8	-	58.9	28.2	-	37.1
CL-34B	55.1	32.1	-	58.2	20.7	-	37.6
CL-34B-Inst	68.7	42.2	-	61.4	27.6	-	40.1
DSC-6.7B	79.1	47.5	-	60.0	30.7	-	38.8
DSC-6.7B-Inst	85.4	50.7	58.1	59.4	32.2	35.6	37.7
DSC-33B	83.4	48.8	-	58.5	32.0	-	38.4
DSC-33B-Inst	88.9	52.1	57.3	58.6	33.8	36.4	38.0
QC-7B	88.0	50.3	57.3	57.2	32.1	34.1	36.5
QC-7B-Inst	89.4	51.3	56.3	57.4	34.2	37.0	38.3
QC-14B	91.0	53.6	58.7	58.9	34.4	35.6	37.8
QC-14B-Inst	90.7	54.3	60.4	59.8	35.8	37.2	39.5
QC-32B	86.7	51.5	59.5	59.3	33.3	37.7	38.4
QC-32B-Inst	89.8	48.2	52.4	53.7	40.7	44.9	45.4

High correctness does not guarantee efficiency. From Table 2, we observe that a high pass rate does not necessarily guarantee superior execution efficiency. For instance, *Claude-4-think* leads in functional correctness with a 95.0% *Pass* score, but it ranks only eighth when considering time efficiency (B_T : 49.8). Notably, it is surpassed by several smaller open-source LLMs, including *QC-14B-Inst*, which achieves the top time efficiency (B_T : 54.3) despite a lower *Pass* score (90.7%). In addition, *CL-34B-Inst*, despite a relatively low pass rate of 68.7%, produces the most time-efficient code on the tasks it translates correctly (B_T^P : 61.4).

To probe the tension between correctness and efficiency, we further conduct controlled correlation analysis. A core challenge here is that efficiency can only be meaningfully measured for functionally correct translations, which inherently couples the two metrics. We therefore adopt two complementary statistical approaches to disentangle this relationship. First, we correlate each model’s *Pass* score with its efficiency on its own set of successful translations (B_T^P/B_M^P). This setup reveals a contradictory trend: correctness is weakly but significantly associated with slower execution time (Kendall’s τ : -0.37, $p < 0.01$), and moderately correlated with higher memory efficiency (Kendall’s τ : 0.52, $p < 0.01$). Second, to eliminate variance arising from models solving different task subsets and establish a standardized comparison, we repeat the analysis on the common set of tasks correctly translated by top-performing models (*Pass* $\geq 85\%$) under consideration (B_T^{com}/B_M^{com}). On this common ground (648 tasks), both correlations shrink and become non-significant (*Time*: τ : -0.28; *Mem.*: τ : 0.28, both

Table 3: Efficiency comparison across code translation directions.

Model	C++→Java		C++→Py		Java→C++		Java→Py		Py→C++		Py→Java		All	
	B_T	B_M												
Claude-4-think	40.9	64.5	56.2	24.8	58.5	51.7	60.8	22.7	49.4	45.5	43.8	62.2	49.8	48.7
Claude-4	38.2	62.4	56.5	20.4	52.3	54.2	56.8	20.1	41.7	41.6	46.1	60.3	47.5	46.7
DS-reasoner	37.1	33.3	40.4	15.8	30.7	23.2	50.0	17.8	32.5	22.4	40.9	26.8	39.0	24.6
DS-chat	39.5	61.0	51.9	18.9	42.2	38.2	60.1	20.9	36.6	32.9	34.9	58.8	43.2	42.9
Gemini-pro	38.2	24.8	40.5	17.9	28.8	22.8	42.6	18.8	25.0	11.7	41.4	27.5	37.3	22.0
Gemini-flash	39.7	31.8	46.3	19.0	29.9	22.9	47.0	16.9	16.4	8.4	34.5	24.5	36.8	22.6
GPT-4o	44.6	58.6	53.0	19.7	34.4	34.7	61.5	18.3	35.5	31.5	41.8	62.1	45.3	42.2
GPT-4o-mini	40.6	59.4	54.3	18.1	44.9	42.7	59.2	15.3	33.6	31.6	42.8	62.0	45.6	42.8
O3	56.4	40.1	58.3	26.7	37.0	25.8	56.8	24.6	25.1	18.0	50.0	37.9	49.4	31.3
O3-mini	43.9	62.6	31.8	45.3	45.7	46.9	53.3	22.9	31.2	32.1	43.4	63.8	42.4	49.5
CL-7b	35.2	26.6	40.5	13.3	24.6	19.3	43.0	13.7	31.4	24.4	37.6	30.7	35.8	22.7
CL-7B-Inst	41.6	34.3	48.6	18.2	43.7	32.8	50.3	19.3	32.2	17.2	46.0	33.3	44.2	27.7
CL-13b	35.6	25.9	39.2	15.7	19.4	18.4	44.3	16.4	18.6	13.0	45.9	30.5	35.9	21.9
CL-13B-Inst	42.8	30.8	48.5	21.1	42.3	34.4	52.2	19.9	32.2	27.2	47.3	31.8	44.8	28.2
CL-34b	37.5	30.4	38.6	15.3	9.5	6.9	48.6	15.3	11.9	9.5	33.8	29.8	32.1	20.7
CL-34B-Inst	39.0	29.0	50.8	25.3	30.7	28.3	53.1	22.8	27.8	21.3	46.6	32.8	42.2	27.6
DSC-6.7B	43.0	34.8	51.7	19.9	45.0	33.9	55.6	19.3	46.7	37.2	46.5	35.0	47.5	30.7
DSC-6.7B-Inst	53.6	39.9	53.1	25.1	42.6	35.2	57.3	20.6	34.1	22.6	54.1	37.9	50.7	32.2
DSC-33B	50.4	43.4	52.5	22.4	35.1	26.8	54.4	18.7	46.4	29.3	50.3	38.2	48.8	32.0
DSC-33B-Inst	57.5	42.8	55.6	25.8	41.0	27.6	59.7	23.7	38.5	32.8	51.8	39.2	52.1	33.8
QC-7B	50.3	37.7	52.8	24.7	48.8	37.6	59.5	15.6	37.1	29.4	49.9	39.0	50.3	32.1
QC-7B-Inst	45.5	38.8	59.4	22.8	47.6	41.0	61.2	27.1	40.3	30.3	53.5	38.5	51.3	34.2
QC-14B	55.0	38.8	59.6	23.5	50.0	46.0	59.1	21.6	42.0	32.2	52.2	38.5	53.6	34.4
QC-14B-Inst	57.2	42.6	57.7	23.6	44.6	38.8	61.3	22.6	39.3	30.4	56.8	44.9	54.3	35.8
QC-32B	54.8	40.8	54.2	22.0	46.8	34.0	55.1	20.6	35.1	29.2	53.9	41.3	51.5	33.3
QC-32B-Inst	59.3	44.6	59.6	22.2	47.9	40.8	27.9	52.9	24.3	31.1	52.9	44.3	48.2	40.7

$p > 0.15$), indicating that the earlier relationships are not robust once task heterogeneity is removed. Taken together, these complementary views suggest that while correctness and efficiency are not completely independent, correctness alone is an unreliable proxy for efficiency. Our findings reveal that efficiency is a distinct capability that should be evaluated independently of functional correctness.

Training strategies yield inconsistent efficiency gains. Common strategies like creating specialized model variants or applying instruction-tuning, while often improving correctness, have an inconsistent impact on efficiency. Among proprietary models, specialized reasoning models do not reliably produce more efficient translation. For instance, while *Claude-4-think* shows a slight efficiency gain over the standard *Claude-4*, the opposite is true for the *DeepSeek* twins. The general-purpose *DS-chat* is substantially more efficient in both time (B_T : 43.2 vs. 39.0) and memory (B_M : 42.9 vs. 24.6) than its specialized *DS-reasoner* counterpart.

Among open-source models, instruction-tuning shows similarly erratic effects. While it consistently improves correctness across model families, its impact on efficiency is not as predictable. For *CL-7B*, instruction-tuning boosts both time and memory efficiency. However, for *QC-32B*, it degrades time efficiency (B_T from 51.5 to 48.2) while improving memory usage (B_M from 33.3 to 40.7), revealing a clear trade-off. Additionally, efficiency does not necessarily scale with model size: the *QC-14B* is more time-efficient than its larger *QC-32B* counterpart (53.6 vs. 51.5). These findings suggest that current tuning strategies, which primarily target correctness, are not yet co-optimized for efficient code translation. This variance also manifests strongly when analyzing performance across different language pairs.

4.2 Granular Efficiency Analysis

To delve deeper, we analyze the model performance across six separate translation directions between C++, Java, and Python. Table 3 summarizes the overall efficiency. A more detailed breakdown is in Appendix A.4 (Table 9).

Efficiency is asymmetric across translation directions.

The choice of source and target language profoundly impacts efficiency, often in asymmetric ways. We observe that models consistently generate more time-efficient code when translating from Java to Python compared to the reverse direction. On average, models score over 7 points higher on overall time efficiency (B_T) for Java→Python (53.5 vs. 46.1). Crucially, this gap is maintained even when controlling for correctness (Table 9, B_T^P : 60.5 vs. 55.0), indicating an intrinsic performance difference in the translated code.

On the other hand, memory efficiency is heavily influenced by the target language’s memory model. From a C++ source, models consistently produce a more memory-efficient translation when targeting Java over Python, with the average B_M being double for C++→Java over C++→Python (41.5 vs. 21.8). The gap persists even after controlling for correctness (Table 9, B_M^P : 49.0 vs. 25.0). Such disparity is rooted in the target languages’ fundamental data representation. Java’s explicit type system and memory-compact primitive types (e.g., *int*) provide a direct mapping from C++, resulting in relatively high memory performance for translated code. In contrast, Python’s *Everything is an object* philosophy introduces a wider spectrum of memory profile possibilities. Consequently, a model’s common literal translation into a memory-heavy Python idiom performs poorly against more optimized alternatives, yielding a low relative score and explaining the lower overall performance.

Table 4: Statistics of classified translation inefficiencies.

Category	Num	%	Time (ΔT)			Mem. (ΔM)		
			Avg	Med.	Max.	Avg	Med.	Max.
Algorithm.	35	14.8	540.1	5.6	6693.7	2.2	1.1	15.0
Idiom./Lib.	146	61.9	2.7	2.1	66.6	3.0	2.1	24.0
Resource.	55	23.3	14.0	2.0	399.1	13.7	12.0	35.4

4.3 Root Case Categorization

To understand the origins of inefficiencies, we first identified a pool of over 2,000 highly inefficient code translations where the performance (*Time* or *Mem.*) is more than twice that of the most efficient for each task. We manually classified a sample of 236 cases (one per task). The resulting taxonomy (Table 10 in Appendix A.5) characterizes inefficiencies into three high-level categories: 1) *Algorithmic-Level Discrepancy*; 2) *Idiomatic and Library-Usage Inefficiency*; 3) *Resource Management and Overhead*. Table 4 summarizes the distribution and impact of these categories.

Our analysis reveals that *Idiomatic and Library-Usage Inefficiency* are the most prevalent, accounting for 61.9% of cases and introducing a median $2.1 \times$ time slowdown (ΔT : 2.1). In stark contrast, *Algorithmic-Level Discrepancies* are the least frequent category (14.8%) but are by far the most destructive, inflicting a catastrophic median time slowdown of $5.6 \times$. Meanwhile, *Resource Management and Overhead* issues (23.3%) are the primary driver of memory inefficiency, increasing median memory consumption by $12.0 \times$. The statistics indicate that idiomatic errors are the most common pitfalls, while the most severe degradation is rooted in categories of algorithm implementation and resource handling, which are less frequent but more detrimental.

Case Study. Figure 3 exemplifies an *Idiomatic/Library-Usage Inefficiency* case from Python→Java translation. The Python code modifies a string using repeated concatenation. On the original benchmark’s small tests (max input string length: 12), two translations from different models demonstrate negligible performance differences (*Time*: 0.164s vs. 0.175s, *Mem.*: 45.0MB vs. 46.0MB). However, under stress tests with larger test inputs (over 10k characters), a massive gap emerges, revealing a $55 \times$ time slowdown and $11 \times$ memory increase in the *Claude-4* version. The root cause is a failure to adopt an idiomatic Java approach. The efficient *Deepseek-chat* translation correctly uses a mutable *StringBuilder* for an $O(L)$ in-place modification. Conversely, *Claude-4* performs a literal translation of Python’s logic. Since Java strings are immutable, such use of repeated concatenation leads to catastrophic $O(L^2)$ time complexity and excessive memory churn.

5 Related Work

Automatic code translation has evolved from early rule-based (Immunant 2025; GoTranspile 2025) and statistical methods (Karaivanov, Raychev, and Vechev 2014; Nguyen, Nguyen, and Nguyen 2016) to neural machine translation (NMT) approaches like the TransCoder series (Lachaux et al. 2020; Roziere et al. 2021; Szafraniec et al. 2022).

```

Source Code In Python
def replace ( s , c1 , c2 ) :
    l = len ( s )
    for i in range ( 1 ) :
        if ( s [ i ] == c1 ) :
            s = s [ 0 : i ] + c2 + s [ i + 1 : ]
        elif ( s [ i ] == c2 ) :
            s = s [ 0 : i ] + c1 + s [ i + 1 : ]
    return s

Java Translation (Deepseek-chat)
static String replace(String s, char c1, char c2) {
    StringBuilder sb = new StringBuilder(s);
    for (int i = 0; i < sb.length(); i++) {
        if (sb.charAt(i) == c1)
            sb.setCharAt(i, c2);
        else if (sb.charAt(i) == c2)
            sb.setCharAt(i, c1);
    }
    return sb.toString();           Execution Time 0.46s
}                                Peak Memory 118.50MB

Java Translation (Claude-4)
static String replace(String s, char c1, char c2) {
    int l = s.length();
    for (int i = 0; i < l; i++) {
        if (s.charAt(i) == c1)
            s = s.substring(0,i)+c2+s.substring(i+1);
        else if (s.charAt(i) == c2)
            s = s.substring(0,i)+c1+s.substring(i+1);
    }
    return s;                      Execution Time 25.26s
}                                Peak Memory 1312.00MB

```

Figure 3: A case study of *Idiom./Lib. Inefficiency*.

The advent of Large Language Models (LLMs) has further advanced the field, with research focusing on improving correctness through techniques like fine-tuning (He et al. 2025; Liu, Li, and Zhang 2023), CoT (Macedo et al. 2024; Nitin, Krishna, and Ray 2024), and self-repair (Yang et al. 2024). Correspondingly, evaluation benchmarks such as Avatar (Ahmad et al. 2021), G-TransEval (Jiao et al. 2023), and PolyHumanEval (Tao et al. 2024) have also prioritized correctness, assessing translation across different language pairs, syntactic complexities, and code granularities. While these efforts have substantially improved correctness, execution efficiency remains a critical but underexplored dimension. This paper introduces the first benchmark to systematically evaluate the efficiency of LLM-based code translation.

6 Conclusion

This work addresses the critical yet overlooked efficiency dimension in LLM-based code translation. We introduce **TRACY**, the first benchmark meticulously designed for its systematic evaluation. Our methodology utilizes an LLM-driven pipeline to generate computationally demanding stress tests and curate efficiency-distinguishing translation tasks, exposing latent performance disparities of up to several orders of magnitude that existing tests would otherwise miss. Our comprehensive evaluation of 26 representative LLMs reveals a stark disconnect between correctness and efficiency: even top models leading in correctness generate correct but inefficient translations that suffer severe performance degradation. To conclude, our findings highlight the urgent need to evolve beyond correctness-centric paradigms and develop novel methods that jointly optimize for both correctness and efficiency in code translation.

References

- Ahmad, W. U.; Chakraborty, S.; Ray, B.; and Chang, K.-W. 2022. Summarize and generate to back-translate: Unsupervised translation of programming languages. *arXiv preprint arXiv:2205.11116*.
- Ahmad, W. U.; Tushar, M. G. R.; Chakraborty, S.; and Chang, K.-W. 2021. Avatar: A parallel corpus for java-python program translation. *arXiv preprint arXiv:2108.11590*.
- Austin, J.; Odena, A.; Nye, M.; Bosma, M.; Michalewski, H.; Dohan, D.; Jiang, E.; Cai, C.; Terry, M.; Le, Q.; et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Chen, M.; Tworek, J.; Jun, H.; Yuan, Q.; Pinto, H. P. D. O.; Kaplan, J.; Edwards, H.; Burda, Y.; Joseph, N.; Brockman, G.; et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Dai, J.; Lu, J.; Feng, Y.; Huang, D.; Zeng, G.; Ruan, R.; Cheng, M.; Tan, H.; and Guo, Z. 2024. Mhpp: Exploring the capabilities and limitations of language models beyond basic code generation. *arXiv preprint arXiv:2405.11430*.
- Du, M.; Luu, A. T.; Ji, B.; Liu, Q.; and Ng, S.-K. 2024. Mercury: A code efficiency benchmark for code large language models. *Advances in Neural Information Processing Systems*, 37: 16601–16622.
- Du, M.; Tuan, L. A.; Liu, Y.; Qing, Y.; Huang, D.; He, X.; Liu, Q.; Ma, Z.; and Ng, S.-k. 2025. Afterburner: Reinforcement Learning Facilitates Self-Improving Code Efficiency Optimization. *arXiv preprint arXiv:2505.23387*.
- Emerson, P. 2013. The original Borda count and partial voting. *Social Choice and Welfare*, 40(2): 353–358.
- GoTranspile. 2025. CxGo: Convert C code to Go. <https://github.com/gotranspile/cxgo>. Accessed: 2025-07-01.
- Guo, D.; Zhu, Q.; Yang, D.; Xie, Z.; Dong, K.; Zhang, W.; Chen, G.; Bi, X.; Wu, Y.; Li, Y.; et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196*.
- He, M.; Yang, F.; Zhao, P.; Yin, W.; Kang, Y.; Lin, Q.; Rajmohan, S.; Zhang, D.; and Zhang, Q. 2025. ExeCoder: Empowering Large Language Models with Executability Representation for Code Translation. *arXiv preprint arXiv:2501.18460*.
- Hendrycks, D.; Basart, S.; Kadavath, S.; Mazeika, M.; Arora, A.; Guo, E.; Burns, C.; Puranik, S.; He, H.; Song, D.; et al. 2021. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*.
- Huang, D.; Dai, J.; Weng, H.; Wu, P.; Qing, Y.; Cui, H.; Guo, Z.; and Zhang, J. 2024a. Effilearn: Enhancing efficiency of generated code via self-optimization. *Advances in Neural Information Processing Systems*, 37: 84482–84522.
- Huang, D.; Qing, Y.; Shang, W.; Cui, H.; and Zhang, J. M. 2024b. Effibench: Benchmarking the efficiency of automatically generated code. *arXiv preprint arXiv:2402.02037*.
- Huang, D.; Zeng, G.; Dai, J.; Luo, M.; Weng, H.; Qing, Y.; Cui, H.; Guo, Z.; and Zhang, J. M. 2024c. SWIFTCODER: Enhancing Code Generation in Large Language Models through Efficiency-Aware Fine-tuning. *arXiv preprint arXiv:2410.10209*.
- Huang, D.; Zhang, J. M.; Bu, Q.; Xie, X.; Chen, J.; and Cui, H. 2024d. Bias testing and mitigation in llm-based code generation. *ACM Transactions on Software Engineering and Methodology*.
- Huang, D.; Zhang, J. M.; Harman, M.; Zhang, Q.; Du, M.; and Ng, S.-K. 2025. Benchmarking LLMs for Unit Test Generation from Real-World Functions. *arXiv preprint arXiv:2508.00408*.
- Hui, B.; Yang, J.; Cui, Z.; Yang, J.; Liu, D.; Zhang, L.; Liu, T.; Zhang, J.; Yu, B.; Lu, K.; et al. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Husein, R. A.; Aburajouh, H.; and Catal, C. 2025. Large language models for code completion: A systematic literature review. *Computer Standards & Interfaces*, 92: 103917.
- Ibrahimzada, A. R.; Ke, K.; Pawagi, M.; Abid, M. S.; Pan, R.; Sinha, S.; and Jabbarvand, R. 2025. AlphaTrans: A Neuro-Symbolic Compositional Approach for Repository-Level Code Translation and Validation. *Proceedings of the ACM on Software Engineering*, 2(FSE): 2454–2476.
- Immunant. 2025. C2Rust: A tool for converting C to Rust. <https://github.com/immunant/c2rust>. Accessed: 2025-07-01.
- Jiang, J.; Wang, F.; Shen, J.; Kim, S.; and Kim, S. 2024. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*.
- Jiao, M.; Yu, T.; Li, X.; Qiu, G.; Gu, X.; and Shen, B. 2023. On the evaluation of neural code translation: Taxonomy and benchmark. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 1529–1541. IEEE.
- Jimenez, C. E.; Yang, J.; Wettig, A.; Yao, S.; Pei, K.; Press, O.; and Narasimhan, K. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Karaivanov, S.; Raychev, V.; and Vechev, M. 2014. Phrase-based statistical translation of programming languages. In *Proceedings of the 2014 ACM international symposium on new ideas, new paradigms, and reflections on programming & software*, 173–184.
- Lachaux, M.-A.; Roziere, B.; Chanussot, L.; and Lample, G. 2020. Unsupervised translation of programming languages. *arXiv preprint arXiv:2006.03511*.
- Li, Y.; Choi, D.; Chung, J.; Kushman, N.; Schrittwieser, J.; Leblond, R.; Eccles, T.; Keeling, J.; Gimeno, F.; Dal Lago, A.; et al. 2022. Competition-level code generation with alphacode. *Science*, 378(6624): 1092–1097.
- Liu, F.; Li, J.; and Zhang, L. 2023. Syntax and domain aware model for unsupervised program translation. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, 755–767. IEEE.
- Liu, J.; Xia, C. S.; Wang, Y.; and Zhang, L. 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: 21558–21572.
- Liu, J.; Xie, S.; Wang, J.; Wei, Y.; Ding, Y.; and Zhang, L. 2024. Evaluating language models for efficient code generation. *arXiv preprint arXiv:2408.06450*.
- Macedo, M.; Tian, Y.; Nie, P.; Cogo, F. R.; and Adams, B. 2024. InterTrans: Leveraging Transitive Intermediate Translations to Enhance LLM-based Code Translation. *arXiv preprint arXiv:2411.01063*.
- Nguyen, T. D.; Nguyen, A. T.; and Nguyen, T. N. 2016. Mapping API elements for code migration with vector representations. In 2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C). IEEE, 756–758.
- Nijkamp, E.; Pang, B.; Hayashi, H.; Tu, L.; Wang, H.; Zhou, Y.; Savarese, S.; and Xiong, C. 2022. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474*.
- Nitin, V.; Krishna, R.; and Ray, B. 2024. Spectra: Enhancing the code translation ability of language models by generating multi-modal specifications. *arXiv preprint arXiv:2405.18574*.
- Niu, C.; Zhang, T.; Li, C.; Luo, B.; and Ng, V. 2024. On evaluating the efficiency of source code generated by llms. In *Proceedings of the 2024 IEEE/ACM First International Conference on AI Foundation Models and Software Engineering*, 103–107.
- Ouyang, S.; Huang, D.; Guo, J.; Sun, Z.; Zhu, Q.; and Zhang, J. M. 2025. DSCodeBench: A Realistic Benchmark for Data Science Code Generation. *arXiv preprint arXiv:2505.15621*.
- Pan, R.; Ibrahimzada, A. R.; Krishna, R.; Sankar, D.; Wassi, L. P.; Merler, M.; Sobolev, B.; Pavuluri, R.; Sinha, S.; and Jabbarvand, R. 2024. Lost in translation: A study of bugs introduced by large language models while translating code. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, 1–13.
- Peng, Y.; Wan, J.; Li, Y.; and Ren, X. 2025. Coffe: A code efficiency benchmark for code generation. *Proceedings of the ACM on Software Engineering*, 2(FSE): 242–265.
- Puri, R.; Kung, D. S.; Janssen, G.; Zhang, W.; Domeniconi, G.; Zolotov, V.; Dolby, J.; Chen, J.; Choudhury, M.; Decker, L.; et al. 2021. Codenet: A large-scale ai for code dataset for learning a diversity of coding tasks. *arXiv preprint arXiv:2105.12655*.
- Qing, Y.; Zhu, B.; Du, M.; Guo, Z.; Zhuo, T. Y.; Zhang, Q.; Zhang, J. M.; Cui, H.; Yiu, S.-M.; Huang, D.; et al. 2025. EffiBench-X: A Multi-Language Benchmark for Measuring Efficiency of LLM-Generated Code. *arXiv preprint arXiv:2505.13004*.
- Roziere, B.; Gehring, J.; Gloeckle, F.; Sootla, S.; Gat, I.; Tan, X. E.; Adi, Y.; Liu, J.; Sauvestre, R.; Remez, T.; et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Roziere, B.; Zhang, J. M.; Charton, F.; Harman, M.; Synnaeve, G.; and Lample, G. 2021. Leveraging automated unit tests for unsupervised code translation. *arXiv preprint arXiv:2110.06773*.
- Seed, B.; Zhang, Y.; Su, J.; Sun, Y.; Xi, C.; Xiao, X.; Zheng, S.; Zhang, A.; Liu, K.; Zan, D.; et al. 2025. Seed-Coder: Let the Code Model Curate Data for Itself. *arXiv preprint arXiv:2506.03524*.
- Shypula, A.; Madaan, A.; Zeng, Y.; Alon, U.; Gardner, J.; Hashemi, M.; Neubig, G.; Ranganathan, P.; Bastani, O.; and Yazdanbakhsh, A. 2023. Learning performance-improving code edits. *arXiv preprint arXiv:2302.07867*.
- Stein’s GCD. 2025. Stein’s Algorithm for Finding GCD. <https://www.geeksforgeeks.org/dsa/steins-algorithm-for-finding-gcd/>. Accessed: 2025-07-01.
- Szafraniec, M.; Roziere, B.; Leather, H.; Charton, F.; Labatut, P.; and Synnaeve, G. 2022. Code translation with compiler representations. *arXiv preprint arXiv:2207.03578*.
- Tao, Q.; Yu, T.; Gu, X.; and Shen, B. 2024. Unraveling the Potential of Large Language Models in Code Translation: How Far are We? In *2024 31st Asia-Pacific Software Engineering Conference (APSEC)*, 353–362. IEEE.
- Vartziotis, T.; Dellatolas, I.; Dasoulas, G.; Schmidt, M.; Schneider, F.; Hoffmann, T.; Kotsopoulos, S.; and Keckisen, M. 2024. Learn to code sustainably: An empirical study on llm-based green code generation. *arXiv preprint arXiv:2403.03344*.
- Xue, P.; Wu, L.; Yang, Z.; Wang, C.; Li, X.; Zhang, Y.; Li, J.; Jin, R.; Pei, Y.; Shen, Z.; et al. 2025. ClassEval-T: Evaluating Large Language Models in Class-Level Code Translation. *Proceedings of the ACM on Software Engineering*, 2(ISSTA): 1421–1444.
- Yang, Z.; Liu, F.; Yu, Z.; Keung, J. W.; Li, J.; Liu, S.; Hong, Y.; Ma, X.; Jin, Z.; and Li, G. 2024. Exploring and unleashing the power of large language models in automated code translation. *Proceedings of the ACM on Software Engineering*, 1(FSE): 1585–1608.
- Zhang, Q.; Fang, C.; Xie, Y.; Ma, Y.; Sun, W.; Yang, Y.; and Chen, Z. 2024. A systematic literature review on large language models for automated program repair. *arXiv preprint arXiv:2405.01466*.
- Zheng, Q.; Xia, X.; Zou, X.; Dong, Y.; Wang, S.; Xue, Y.; Shen, L.; Wang, Z.; Wang, A.; Li, Y.; et al. 2023a. 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.
- Zheng, Z.; Ning, K.; Wang, Y.; Zhang, J.; Zheng, D.; Ye, M.; and Chen, J. 2023b. A survey of large language models for code: Evolution, benchmarking, and future trends. *arXiv preprint arXiv:2311.10372*.
- Zhu, M.; Karim, M.; Lourentzou, I.; and Yao, D. 2024. Semi-supervised code translation overcoming the scarcity of parallel code data. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, 1545–1556.

A Appendix

A.1 Details of Stress Test Generation

Table 5 presents the average execution time and peak memory usage for the top-10 stress tests identified across five iterations. Both metrics consistently increase with each iteration, confirming that our iterative efficiency-oriented mechanism successfully guides the LLM to discover more resource-intensive tests. We set the maximum number of iterations to five, as further iterations began to yield tests so computationally demanding that they caused some ground-truth implementations from the Transcoder-Test to exceed time or memory limits, indicating diminishing returns.

Table 5: The average execution time and peak memory usage of the top-10 stress tests after each iteration.

Language		Iter 1	Iter 2	Iter 3	Iter 4	Iter 5
C++	Time	0.136	0.235	0.334	0.423	0.504
	Mem	19.010	22.432	25.854	28.695	31.663
Java	Time	0.495	0.687	0.879	1.082	1.270
	Mem	131.098	173.689	216.280	252.236	279.411
Python	Time	1.447	3.928	6.409	8.287	9.883
	Mem	54.128	584.995	1115.862	1496.826	1897.970
Avg.	Time	0.675	1.588	2.501	3.202	3.803
	Mem	67.355	256.954	446.553	582.274	721.087

Table 6 presents a direct comparison of the execution profiles between the original tests and our generated stress tests on the most-efficient translations. As shown, the generated tests significantly amplify the computational demands. For the initial pool of tasks before pruning, our tests increase the average execution time by a factor of $48.8\times$ (from 0.078s to 3.803s) and peak memory usage by $42.3\times$ (from 17.05MB to 721.087MB). Even after pruning, the stress tests reveal a significant performance gap on the most efficient translations, increasing execution time by $9.4\times$ and memory usage by $4.0\times$. This substantial increase in consumption is critical for exposing latent efficiency differences that are undetectable with the original, less demanding tests.

Table 6: Execution comparison of top-10 stress tests with original tests on the most-efficient reference translations.

Tests	C++		Java		Python		Average	
	Time	Mem.	Time	Mem.	Time	Mem.	Time	Mem.
<i>Before Pruning</i>								
Orig.	0.002	1.506	0.209	43.609	0.018	7.368	0.078	17.793
Stress	0.506	31.507	1.270	279.411	9.916	1897.970	3.803	721.087
<i>After Pruning</i>								
Orig.	0.001	1.507	0.210	43.609	0.010	5.458	0.094	22.008
Stress	0.447	25.219	0.610	159.006	2.851	111.791	0.883	89.010

The effectiveness of these generated stress tests is further validated in our root-cause analysis of 236 inefficient translation cases. We confirmed that the original benchmark’s simpler tests overlooked 82.6% of the performance degra-

dation issues that TRACY’s more computationally demanding tests identified. This significant disparity highlights that stress tests are essential for a genuine assessment of code translation efficiency, as they successfully expose severe performance bottlenecks that would otherwise go unnoticed.

A.2 Details of Task Pruning

Table 7 details the task pruning pipeline, showing the number of tasks filtered at each stage to curate the final TRACY benchmark. Starting with an initial pool of 2,828 translation tasks, we first discard 362 tasks for which a valid stress test suite (\mathcal{T}_s) could not be generated. Subsequently, we apply the three filtering criteria outlined in the main paper: Rule-1 (Feasibility) removes 196 tasks that no model could solve correctly; Rule-2 (Impactfulness) discards 311 tasks where efficiency was not a practical concern; and Rule-3 (Diversity) filters out 948 tasks where all correct solutions had negligible performance differences. This multi-stage filtering process ensures that the resulting 1,011 tasks in TRACY are genuinely efficiency-critical.

Table 7: Details of filtered translation tasks by each rule.

Translation	Orig.	w/o \mathcal{T}_s	Rule-1	Rule-2	Rule-3	TRACY
cpp→Java	482	59	12	0	161	250
cpp→Py	464	63	30	0	231	140
Java→cpp	468	59	19	160	97	133
Java→Py	464	63	27	0	227	147
Py→cpp	468	59	75	151	79	104
Py→Java	482	59	33	0	153	237
All	2828	362	196	311	948	1011

A.3 Model List

Table 8 provides details for the 26 Large Language Models (LLMs) used in our work. These models were utilized for both the curation of the benchmark and the final evaluation. During the efficiency-oriented task pruning stage, we generated a diverse set of candidate translations by sampling one translation from each model for each task, using a temperature of 0.8. For the final evaluation, we generate one translation from each model via greedy decoding with zero-shot prompting to assess the default translation performance.

A.4 Extended Results

We provide a granular breakdown of model performance to supplement the findings in Section 4.2. Table 9 details the results for each of the six translation directions to support our analysis of asymmetric and language-dependent efficiency. The table reports the correctness (*Pass*), overall efficiency scores (B_T/B_M), and the efficiency scores conditioned on correctness (B_T^P/B_M^P), offering a comprehensive view of model capabilities across different language pairs.

A.5 Taxonomy

To systematically classify the root causes of performance issues, we developed a taxonomy of inefficiency patterns.

This taxonomy is grounded in a rigorous manual analysis of 236 inefficient translations (one per task among all identified inefficiencies). The process was conducted by two of the authors in this work, who independently annotated each case. Initial disagreements were resolved through discussion until a consensus was reached for every case. Through this iterative process of labeling and refinement, we established the three high-level categories presented in Table 10. These categories, which align with our main quantitative findings, provide a clear framework for understanding performance degradation in LLM-based code translation, ranging from high-level algorithmic flaws to low-level resource management errors.

A.6 Extended Discussion on Related Work

Code LLMs Recent years have witnessed a prosperous advancement in Large Language Models (LLMs) tailored for coding tasks. Models such as Codex (Chen et al. 2021), AlphaCode (Li et al. 2022) and more recently, specialized models like CodeLlama (Roziere et al. 2023), DeepSeek-Coder (Guo et al. 2024), Seed-Coder (Seed et al. 2025) and Qwen2.5-Coder (Hui et al. 2024) have demonstrated remarkable capabilities in code generation (Jiang et al. 2024) and debugging (Zhang et al. 2024). These models leverage vast code corpora to learn intricate programming patterns and semantic relationships, significantly revolutionizing software development by automating complex coding tasks.

Benchmarks in Code Generation and Translation The evaluation of code LLMs has been driven by numerous influential benchmarks (Chen et al. 2021; Austin et al. 2021; Hendrycks et al. 2021; Liu et al. 2023; Qing et al. 2025; Huang et al. 2024b; Ouyang et al. 2025; Huang et al. 2024d; Dai et al. 2024; Huang et al. 2025). For code generation, foundational benchmarks include HumanEval (Chen et al. 2021) and MBPP (Austin et al. 2021), which introduce entry-level programming puzzles. To assess more complex reasoning, APPS (Hendrycks et al. 2021) curated thousands of problems from competitive programming platforms. Subsequent efforts like HumanEval+ (Liu et al. 2023) and MHPP (Dai et al. 2024) further augmented these benchmarks with more comprehensive tests to better detect subtle bugs.

For code translation, benchmarks have traditionally prioritized correctness over efficiency. The TransCoder-Test (Lachaux et al. 2020; Yang et al. 2024) benchmark, on which our work builds, established a standard by providing parallel functions in C++, Java, and Python. More targeted benchmarks like Avatar (Ahmad et al. 2022) focused on Java-Python translation by mining parallel methods from open-source projects. To enable a more fine-grained analysis, G-TransEval (Jiao et al. 2023) categorized translation tasks based on four levels of syntactic complexity. Poly-HumanEval (Tao et al. 2024) offers more linguistic diversity by extending HumanEval to 14 programming languages. More recently, ClassEval-T (Xue et al. 2025) and AlphaTrans (Ibrahimzada et al. 2025) have advanced the field towards higher complexity by benchmarking LLM un-

der class- and repository-level code translation. However, a common limitation across these benchmarks is the lack of tests designed to stress the performance of translated code, leaving the efficiency dimension unexplored.

Efficiency of LLM-generated Code While early evaluations of LLM-generated code focused on functional correctness (Chen et al. 2021; Liu et al. 2024), the crucial dimension of execution efficiency has recently gained significant attention in code generation. Several benchmarks have been developed to address this gap by creating challenging testbeds and novel metrics. Initiatives like EffiBench (Huang et al. 2024b) and EffiBench-X (Qing et al. 2025) established a systematic framework by curating efficiency-critical problems with highly optimized, human-written reference solutions. Similarly, the EvalPerf benchmark (Liu et al. 2024) utilized a Differential Performance Evaluation (DPE) framework, which generates computationally expensive inputs to stress solutions and evaluates their performance against a spectrum of reference implementations with varying efficiency levels. Other efforts have focused on developing new evaluation metrics. Mercury (Du et al. 2024) introduced the *Beyond* metric, which integrates correctness with a runtime percentile score derived from a large pool of solutions.

Beyond evaluation, a parallel line of research has focused on directly optimizing the efficiency of LLM-generated code. For instance, SOAP (Huang et al. 2024a) introduced a self-optimization framework where an LLM iteratively improves the code efficiency from performance feedback. Other approaches focus on teaching models to generate efficient code from the outset; PIE (Shypula et al. 2023) and SwiftCoder (Huang et al. 2024c) achieve this by fine-tuning models on datasets of code paired with more efficient solutions. More recently, Afterburner (Du et al. 2025) leveraged reinforcement learning to train an agent that applies code transformations to improve runtime performance.

Our work is situated at the intersection of these research areas but is distinct in its focus. While recent work has begun to address the efficiency of code generation, the unique challenges of efficiency in code translation, where the model must interpret, preserve, and adapt idiomatic properties from a source language, remain unexplored. We introduced TRACY, the first benchmark designed specifically to evaluate this critical translation-specific efficiency dimension.

Table 8: Model details, including aliases used in this paper, their original names, parameter scales, and public links.

Alias In the Paper	Model Name	Model Scale	Link
<i>Proprietary Models</i>			
Claude-4-think	Claude 4 (Thinking)	–	https://www.anthropic.com/news/clause-4
Claude-4	Claude 4 (Non-thinking)	–	https://www.anthropic.com/news/clause-4
DS-reasoner	DeepSeek Reasoner	–	https://www.deepseek.com
DS-chat	DeepSeek Chat	–	https://www.deepseek.com
Gemini-pro	Gemini 2.5 Pro	–	https://deepmind.google/models/gemini/pro
Gemini-flash	Gemini 2.5 Flash	–	https://deepmind.google/models/gemini/flash
GPT-4o	GPT-4o	–	https://openai.com
GPT-4o-mini	GPT-4o-mini	–	https://openai.com
O3	OpenAI o3	–	https://openai.com
O3-mini	OpenAI o3-mini	–	https://openai.com
<i>Open-Source Models</i>			
CL-7B	CodeLlama-7B	7B	https://github.com/facebookresearch/codellama
CL-7B-Inst	CodeLlama-7B-Instruct	7B	https://github.com/facebookresearch/codellama
CL-13B	CodeLlama-13B	13B	https://github.com/facebookresearch/codellama
CL-13B-Inst	CodeLlama-13B-Instruct	13B	https://github.com/facebookresearch/codellama
CL-34B	CodeLlama-34B	34B	https://github.com/facebookresearch/codellama
CL-34B-Inst	CodeLlama-34B-Instruct	34B	https://github.com/facebookresearch/codellama
DSC-6.7B	deepseek-coder-6.7b	6.7B	https://github.com/deepseek-ai/DeepSeek-Coder
DSC-6.7B-Inst	deepseek-coder-6.7b-Instruct	6.7B	https://github.com/deepseek-ai/DeepSeek-Coder
DSC-33B	deepseek-coder-33b	33B	https://github.com/deepseek-ai/DeepSeek-Coder
DSC-33B-Inst	deepseek-coder-33b-Instruct	33B	https://github.com/deepseek-ai/DeepSeek-Coder
QC-7B	Qwen2.5-Coder-7B	7B	https://github.com/QwenLM/Qwen2.5-Coder
QC-7B-Inst	Qwen2.5-Coder-7B-Instruct	7B	https://github.com/QwenLM/Qwen2.5-Coder
QC-14B	Qwen2.5-Coder-14B	14B	https://github.com/QwenLM/Qwen2.5-Coder
QC-14B-Inst	Qwen2.5-Coder-14B-Instruct	14B	https://github.com/QwenLM/Qwen2.5-Coder
QC-32B	Qwen2.5-Coder-32B	32B	https://github.com/QwenLM/Qwen2.5-Coder
QC-32B-Inst	Qwen2.5-Coder-32B-Instruct	32B	https://github.com/QwenLM/Qwen2.5-Coder

Table 9: Granular correctness and efficiency performance across six translation directions. For each direction, we report the Pass rate (Pass), overall and conditional *Beyond* for time (B_T, B_T^P), and overall and conditional *Beyond* for memory (B_M, B_M^P).

Model	C++→Java					C++→Py					Java→C++				
	Pass	B_T	B_T^P	B_M	B_M^P	Pass	B_T	B_T^P	B_M	B_M^P	Pass	B_T	B_T^P	B_M	B_M^P
<i>Proprietary Models</i>															
Claude-4-think	95.2	40.9	42.9	64.5	67.7	97.1	56.2	57.9	24.8	25.5	97.7	58.5	59.9	51.7	52.9
Claude-4	96.8	38.2	39.5	62.4	64.5	97.1	56.5	58.2	20.4	21.0	90.2	52.3	58.0	54.2	60.0
DS-reasoner	80.8	37.1	46.0	33.3	41.2	77.9	40.4	51.9	15.8	20.2	60.2	30.7	51.1	23.2	38.5
DS-chat	96.8	39.5	40.8	61.0	63.0	95.7	51.9	54.2	18.9	19.7	81.2	42.2	52.0	38.2	47.1
Gemini-pro	58.8	38.2	65.0	24.8	42.1	70.7	40.5	57.3	17.9	25.3	55.6	28.8	51.7	22.8	41.0
Gemini-flash	68.4	39.7	58.1	31.8	46.5	73.6	46.3	62.9	19.0	25.9	52.6	29.9	56.8	22.9	43.5
GPT-4o	94.8	44.6	47.1	58.6	61.8	97.1	53.0	54.5	19.7	20.2	64.7	34.4	53.2	34.7	53.7
GPT-4o-mini	90.8	40.6	44.8	59.4	65.4	95.0	54.3	57.2	18.1	19.1	78.2	44.9	57.5	42.7	54.6
O3	93.6	56.4	60.2	40.1	42.8	95.0	58.3	61.3	26.7	28.1	65.4	37.0	56.6	25.8	39.4
O3-mini	94.8	43.9	46.4	62.6	66.0	95.7	31.8	33.2	45.3	47.3	85.7	45.7	53.3	46.9	54.7
<i>Open-Source Models</i>															
CL-7B	56.8	35.2	61.9	26.6	46.8	67.1	40.5	60.4	13.3	19.8	45.9	24.6	53.6	19.3	42.0
CL-7B-Inst	74.8	41.6	55.6	34.3	45.8	82.1	48.6	59.2	18.2	22.1	69.9	43.7	62.5	32.8	46.9
CL-13B	62.4	35.6	57.0	25.9	41.5	67.9	39.2	57.8	15.7	23.1	36.1	19.4	53.8	18.4	50.9
CL-13B-Inst	72.8	42.8	58.8	30.8	42.3	80.7	48.5	60.0	21.1	26.2	75.9	42.3	55.7	34.4	45.3
CL-34B	66.4	37.5	56.5	30.4	45.8	62.1	38.6	62.1	15.3	24.6	14.3	9.5	66.4	6.9	48.0
CL-34B-Inst	65.6	39.0	59.4	29.0	44.1	82.9	50.8	61.3	25.3	30.6	49.6	30.7	61.9	28.3	57.1
DSC-6.7B	76.4	43.0	56.3	34.8	45.6	87.9	51.7	58.8	19.9	22.7	74.4	45.0	60.4	33.9	45.5
DSC-6.7B-Inst	90.0	53.6	59.6	39.9	44.3	89.3	53.1	59.4	25.1	28.1	72.9	42.6	58.4	35.2	48.3
DSC-33B	88.8	50.4	56.7	43.4	48.9	87.1	52.5	60.3	22.4	25.7	66.9	35.1	52.4	26.8	40.0
DSC-33B-Inst	94.8	57.5	60.7	42.8	45.1	93.6	55.6	59.4	25.8	27.5	74.4	41.0	55.1	27.6	37.1
QC-7B	89.2	50.3	56.4	37.7	42.3	92.9	52.8	56.8	24.7	26.6	82.7	48.8	59.0	37.6	45.4
QC-7B-Inst	93.6	45.5	48.7	38.8	41.5	92.1	59.4	64.4	22.8	24.7	84.2	47.6	56.5	41.0	48.6
QC-14B	91.2	55.0	60.3	38.8	42.5	94.3	59.6	63.3	23.5	24.9	87.2	50.0	57.4	46.0	52.8
QC-14B-Inst	95.6	57.2	59.8	42.6	44.5	97.1	57.7	59.4	23.6	24.3	80.5	44.6	55.4	38.8	48.3
QC-32B	91.6	54.8	59.9	40.8	44.5	91.4	54.2	59.2	22.0	24.0	82.0	46.8	57.1	34.0	41.5
QC-32B-Inst	94.8	59.3	62.6	44.6	47.0	93.6	59.6	63.6	22.2	23.8	82.0	47.9	58.5	40.8	49.8
Model	Java→Py					Py→C++					Py→Java				
	Pass	B_T	B_T^P	B_M	B_M^P	Pass	B_T	B_T^P	B_M	B_M^P	Pass	B_T	B_T^P	B_M	B_M^P
<i>Proprietary Models</i>															
Claude-4-think	99.3	60.8	61.2	22.7	22.8	82.7	49.4	59.7	45.5	55.1	94.5	43.8	46.3	62.2	65.8
Claude-4	98.0	56.8	58.0	20.1	20.5	77.9	41.7	53.5	41.6	53.4	97.0	46.1	47.5	60.3	62.1
DS-reasoner	80.3	50.0	62.3	17.8	22.2	55.8	32.5	58.3	22.4	40.2	68.4	40.9	59.8	26.8	39.1
DS-chat	92.5	60.1	65.0	20.9	22.6	66.3	36.6	55.2	32.9	49.5	90.7	34.9	38.5	58.8	64.8
Gemini-pro	73.5	42.6	58.0	18.8	25.6	40.4	25.0	62.0	11.7	29.1	66.2	41.4	62.5	27.5	41.4
Gemini-flash	78.9	47.0	59.6	16.9	21.4	27.9	16.4	58.8	8.4	30.3	56.5	34.5	61.1	24.5	43.3
GPT-4o	97.3	61.5	63.2	18.3	18.8	60.6	35.5	58.6	31.5	52.0	96.6	41.8	43.2	62.1	64.2
GPT-4o-mini	97.3	59.2	60.9	15.3	15.8	63.5	33.6	52.9	31.6	49.8	94.1	42.8	45.5	62.0	65.9
O3	94.6	56.8	60.1	24.6	26.0	49.0	25.1	51.3	18.0	36.8	86.9	50.0	57.5	37.9	43.6
O3-mini	95.2	53.3	55.9	22.9	24.0	58.7	31.2	53.2	32.1	54.7	94.1	43.4	46.1	63.8	67.8
<i>Open-Source Models</i>															
CL-7B	70.7	43.0	60.8	13.7	19.4	53.8	31.4	58.4	24.4	45.3	68.8	37.6	54.6	30.7	44.7
CL-7B-Inst	81.6	50.3	61.6	19.3	23.6	51.0	32.2	63.2	17.2	33.7	79.3	46.0	58.0	33.3	42.0
CL-13B	70.7	44.3	62.7	16.4	23.2	31.7	18.6	58.6	13.0	41.0	72.2	45.9	63.7	30.5	42.2
CL-13B-Inst	85.0	52.2	61.4	19.9	23.4	56.7	32.2	56.8	27.2	48.0	80.2	47.3	59.0	31.8	39.7
CL-34B	76.2	48.6	63.8	15.3	20.1	18.3	11.9	64.9	9.5	52.0	65.0	33.8	51.9	29.8	45.9
CL-34B-Inst	83.0	53.1	64.0	22.8	27.4	44.2	27.8	62.9	21.3	48.1	76.4	46.6	61.0	32.8	42.9
DSC-6.7B	85.7	55.6	64.9	19.3	22.5	71.2	46.7	65.7	37.2	52.3	78.9	46.5	58.9	35.0	44.3
DSC-6.7B-Inst	93.9	57.3	61.0	20.6	22.0	56.7	34.1	60.1	22.6	39.8	92.4	54.1	58.5	37.9	41.0
DSC-33B	87.1	54.4	62.5	18.7	21.4	72.1	46.4	64.3	29.3	40.7	87.3	50.3	57.5	38.2	43.7
DSC-33B-Inst	96.6	59.7	61.8	23.7	24.6	71.2	38.5	54.1	32.8	46.1	91.1	51.8	56.9	39.2	43.0
QC-7B	93.9	59.5	63.4	15.6	16.6	66.3	37.1	55.8	29.4	44.3	92.8	49.9	53.7	39.0	42.0
QC-7B-Inst	92.5	61.2	66.2	27.1	29.3	72.1	40.3	55.9	30.3	42.0	92.0	53.5	58.2	38.5	41.9
QC-14B	95.9	59.1	61.6	21.6	22.5	75.0	42.0	56.0	32.2	42.9	94.9	52.2	55.0	38.5	40.6
QC-14B-Inst	97.3	61.3	63.0	22.6	23.3	66.3	39.3	59.3	30.4	45.8	94.1	56.8	60.4	44.9	47.7
QC-32B	90.5	55.1	60.9	20.6	22.8	59.6	35.1	58.9	29.2	48.9	91.1	53.9	59.1	41.3	45.3
QC-32B-Inst	96.6	27.9	28.9	52.9	54.7	60.6	24.3	40.1	31.1	51.4	95.4	52.9	55.5	44.3	46.4

Table 10: A Taxonomy of inefficiencies in LLM-based code translation.

1. Algorithmic-Level Discrepancy
1.1 Asymptotic Complexity Degradation <i>Description.</i> Translates an efficient algorithm to a higher-complexity variant. <i>Example:</i> Translating an $O(1)$ hash table lookup into an $O(N^3)$ triple-nested loop.
1.2 Failure to Adopt Idiomatic Algorithm <i>Description.</i> Preserves the source algorithm when a more efficient idiomatic alternative exists in the target language. <i>Example:</i> Literally translating an $O(N \log N)$ Python sort to find the maximum instead of using C++ <code>std::max_element(O(N))</code> .
1.3 Flawed Semantic Refactoring <i>Description.</i> Introduces logical errors when attempting to refactor code to fit the target language during translation. <i>Example:</i> Translating C++ integer division <code>a/b</code> to Python's <code>a/b</code> (float division). Such internal semantic alteration sometimes still leads to the correct output, but may incur an efficiency penalty from unnecessary floating-point operations within loops.
2. Idiomatic & Library-Usage Inefficiency
2.1 Suboptimal Data Structure Choice <i>Description.</i> Chooses a data structure that is functionally correct but asymptotically or practically slower in the target language. <i>Example:</i> Translating Java's hash map <code>HashMap</code> to C++ <code>std::map(O(\log N))</code> instead of <code>std::unordered_map</code> (expected $O(1)$).
2.2 Inefficient API / Function Usage <i>Description.</i> Utilizes high-overhead APIs or patterns, particularly within efficiency-critical hot loops or large-scale data flows. <i>Example:</i> Building strings via <code>+= concatenation(O(N^2))</code> instead of using Java's mutable <code>StringBuilder</code> or Python's <code>''.join()</code> .
2.3 Misunderstanding of Core Semantics <i>Description.</i> Misinterprets value vs. reference semantics, copying, or mutability, causing excessive allocations or unintended side effects. <i>Example:</i> Initializing a 2D Python list via <code>[[0]*n]*n</code> , which creates shallow copies of inner lists and causes aliasing.
3. Resource Management and Overhead
3.1 Memory Allocation Mismatch <i>Description.</i> Employs allocation strategies that are non-standard for the target language, causing heap churn or fragmentation. <i>Example:</i> Using <code>new/delete</code> arrays instead of <code>std::vector</code> in C++, leading to fragmented allocations and manual lifetime costs.
3.2 Unnecessary Abstraction Overhead <i>Description.</i> Introduces heavyweight abstractions for simple operations, incurring unnecessary runtime or memory costs. <i>Example:</i> Using Java's <code>BigInteger</code> for arithmetic that fits within a primitive <code>long</code> , causing heap allocations in tight loops.
3.3 Inefficient Type Mapping <i>Description.</i> Translates compact data types into less memory-efficient representations in the target language. <i>Example:</i> Mapping C++ <code>std::vector<bool></code> to Java <code>Vector<Boolean></code> , triggering per-element object boxing and GC pressure.