

LEGO-COMPILER: ENHANCING NEURAL COMPILATION THROUGH COMPOSABLE CHAIN OF THOUGHT

Anonymous authors
Paper under double-blind review

ABSTRACT

Large language models (LLMs) have the potential to revolutionize how we design and implement compilers and code translation tools. However, existing LLMs struggle to handle long and complex programs. We introduce LEGO-Compiler, a novel neural compilation system that leverages LLMs to translate high-level languages into assembly code. Our approach centers on three key innovations: LEGO translation, which decomposes the input program into manageable blocks; annotation-based Chain-of-Thoughts, guiding LLMs through the compilation process with LLM-annotated context; and a feedback mechanism for self-correction. Supported by formal proofs of code composability, LEGO-Compiler demonstrates high accuracy on multiple datasets, including over 99% on ExeBench and 100% on industrial-grade CoreMark, and successfully handles programs far exceeding the length limitations of native LLM translation. This work opens new avenues for applying LLMs to system-level tasks, complementing traditional compiler technologies.

1 INTRODUCTION

The rapid development of Large Language Models (LLMs) has led to an expansion of their applications and effectiveness across various domains (Rombach et al., 2022; OpenAI, 2023; 2024; Ziegler et al., 2024). One important area where LLMs have shown impressive results is code translation, including tasks such as code generation from natural languages (Zan et al., 2023) and transformation between programming languages (Yang et al., 2024). In code translation, LLMs have demonstrated remarkable accuracy and readability, often surpassing manually crafted translators.

While LLMs have shown promising results in translating between high-level programming languages (Rozière et al., 2020; Roziere et al., 2021; Szafraniec et al., 2023) and in decompilation tasks (Fu et al., 2019; Cao et al., 2022; Armengol-Estapé et al., 2023), their application to translating from high-level languages to low-level assembly languages remains a relatively unexplored area. This can be attributed to two main factors. Firstly, the dominance of traditional compilers in this domain has left little incentive for exploring alternatives in such a mature field. Secondly, the complexity of compilers and the bitwise precision required in compilation tasks for semantic accuracy have made it challenging for LLMs based on statistical learning.

Despite these challenges, LLMs have shown promising capabilities in compilation-related tasks. Cummins et al. (2023; 2024) has demonstrated their proficiency in optimizing compiler options and their excellent ability to mimic compiler code behavior, producing high-quality IR code. Furthermore, preliminary explorations in translating from high-level languages to assembly languages, such as C-x86 (Armengol-Estapé & O’Boyle, 2021) and C-LLVM IR (Guo & Moses, 2022), have indicated the potential feasibility of using LLMs in compilation tasks. However, these existing works have not fully addressed the boundaries of LLM capabilities in compilation tasks – specifically, what LLMs can and cannot do in this domain. Compilation is typically divided into two main aspects: functionality and optimization. This work focuses on exploring and answering questions about LLM capabilities in the functionality aspect of compilation.

LLMs are pre-trained on vastly large code corpora. some are monolingual, and some may be bilingual (where LLMs can learn the translation rule between two languages). However, most of these LLMs do not disclose their training datasets, so their capabilities can only be assessed through

empirical testing. We primarily find that LLMs learn the neural compilation process from directly compiler-generated bilingual corpora, which is a relatively easy way to construct pretraining dataset. However, we found that assembly code directly generated by traditional compilers is hard to learn for LLM-based generation due to several challenges. These include the presence of semantically opaque labels, symbols or numeric values that LLMs struggle to translate accurately, and the need to handle symbol renaming for identifiers with the same name in different scopes, etc. Although style migration or modifications to existing compilers can be made, these approaches still rely on an existing compiler to perform the neural compilation job, which doesn't outperform existing designs.

Our work takes a different approach where we do not require bilingual corpora. As a result, we don't rely on an existing compiler. Regarding LLMs' strong in-context learning abilities (Min et al., 2022; Song et al., 2024), we propose the following methods: through high-quality examples and compiler knowledge guided Chain-of-Thoughts, LLMs can perform step-by-step neural compilation. This approach involves generating annotations highly corresponding to source code statements and data structure layout annotations, leading to substantial improvements in the compilation generation task.

More importantly, the scalability of current code translation is also a big problem. Although advanced LLMs already have hundreds of thousands tokens context limit, they can not merely compile a code with 2.6k tokens in CoreMark (Gal-On & Levy, 2012), which is just a 200-LOC function. To address this limitation, we have an intuitive thought: can we split the program into finer-grained components, compile each component, then assemble them together? We can surely do it in function-level, since functions are trivial semantic units (Ibrahimzada, 2024). However, we still think function-level is still too coarse-grained, and we seek to further break it.

Based on these insights, we propose a novel approach called LEGO translation, which draws inspiration from the modular and composable nature of LEGO blocks. This method breaks down large programs into manageable, semantically-composable control blocks, analogous to LEGO pieces. These blocks are then independently translated and rebuilt to form a much larger scale translation. We apply the LEGO translation method to the compilation domain and, guided by a series of compilation-specific expert knowledge, design the LEGO-Compiler, a scalable, LLM-driven system that leverages the power of LLMs to perform neural compilation tasks.

LEGO-Compiler can correctly compile over 99% of the code in ExeBench (Armengol-Estapé et al., 2022), a large scale dataset through careful unit-testing. We can also correctly compile 100% CoreMark (Gal-On & Levy, 2012), an industrial-grade code that encompasses most common programming language features in C. Regarding scalability, we have verified that LEGO translation method can significantly scale up the capability of code translation/compilation performed by LLMs, where we propose the LongFunction dataset for very long code translation and compilation evaluation, and the LEGO translation method can sufficiently translate it for both neural compilation or code translation purposes.

The main contributions of this work are as follows:

- We propose a set of novel methods for neural compilation, including LEGO translation for breaking down large programs, annotation-based Chain-of-Thoughts (CoTs) that explicitly generate intermediate results to aid translation, and a feedback-driven self-correction mechanism. These methods collectively address the challenges of applying LLMs to complex compilation tasks.
- We introduce LEGO-Compiler, a comprehensive neural compilation system that integrates our proposed methods. LEGO-Compiler incorporates control flow annotation, struct annotation, and variable mapping to ensure accurate and scalable compilation across various architectures and programming languages.
- We provide both theoretical and empirical support for our approach. We present a formal proof of code composability that underpins the LEGO translation method. Empirically, we demonstrate LEGO-Compiler's effectiveness through extensive evaluations, achieving over 99% accuracy on ExeBench and 100% accuracy on the industrial-grade CoreMark benchmark. Our system successfully handles programs far exceeding the length limitations of direct LLM translation, showcasing its scalability.

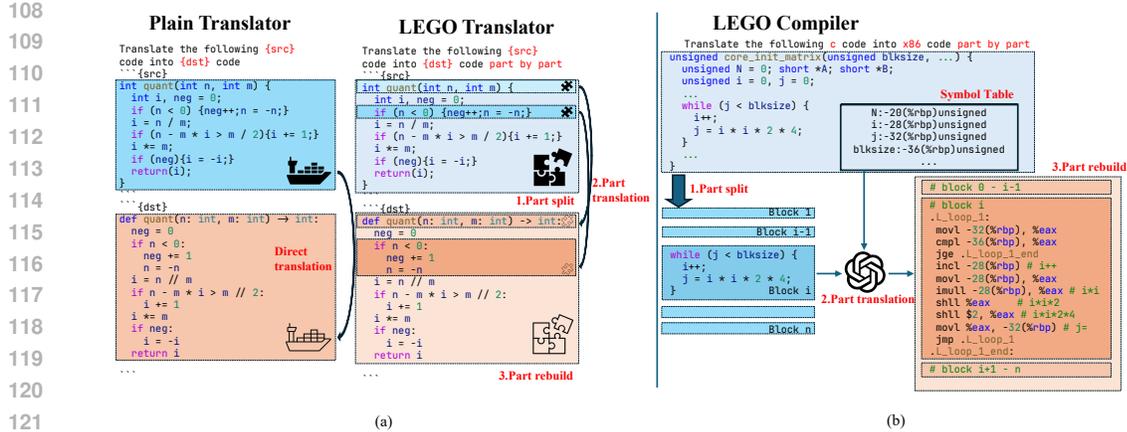


Figure 1: **a.** Plain translation vs LEGO translation, by splitting the program into smaller composable control blocks(parts), translating each part becomes an easier task, and rebuilding each translated partial result will form a full translation. **b.** LEGO compiler, a special case for LEGO translation, to translate each part correctly, a symbol table need to be maintained first and provided during translation.

2 METHODS

2.1 PROBLEM DEFINITION

Before introducing our method, we first define the neural compilation problem. Neural compilation can be viewed as a specialized version of code translation problem, as defined in Definition 1, with the goal of translating high-level programming language as the src language (such as C) into low-level assembly language as the dst language (such as x86, ARM, or RISC-V). Unlike general code translation, compilation needs to handle more low-level details, such as memory layout and calling convention, while ensuring the functional correctness of the translated result.

Definition 1. There are two programming languages: \mathcal{L}_{src} and \mathcal{L}_{dst} , each is an infinite set of valid program strings. There exists a unary relation \rightarrow from \mathcal{L}_{src} to \mathcal{L}_{dst} . The problem is to perform a translator function $T: \forall x \in \mathcal{L}_{src}, (\exists u \in \mathcal{L}_{dst}, x \rightarrow u) \rightarrow (x \rightarrow T(x)), T(x) \equiv x$ semantically.

2.2 LEGO TRANSLATION: CORE METHOD

As depicted in (a) in Figure 1, previous neural code translation methods typically convert entire programs at the function or file level. While this approach may be effective for smaller programs, it struggles with larger programs due to significant accuracy degradation. These methods translate code at a coarse granularity, making it challenging to translate very long functions using LLMs. Taking neural compilation as an example, all current LLMs fail to compile a C function with larger than 2.6k tokens using direct translation, although some advanced LLMs already have 128k-200k context limit. They could also perform code-snippet level translation, but they lack guidelines and necessary information to compose the code-snippet level results together, and there is also no clear formal proof to the composability of code. Despite these limitations, we observe an inherently composable nature in code. In the context of neural compilation, we propose the following insights to enhance translation scalability:

- **Fine-grained translation:** Instead of translating an entire program at once, focus on translating smaller code snippets accurately. By ensuring each part is correctly translated, they can be combined to form a semantically equivalent complete translation.
- **Contextual Awareness:** Effective translation of smaller code snippets requires understanding their contextual positioning within the code. This includes recognizing the relationship with preceding and succeeding snippets to maintain semantic coherence.

- **Symbol Handling:** Accurate translation involves reasoning about necessary symbols and constructs within each snippet to enhance alignment with the target language’s syntax and semantics, thereby aiming to preserve the intended functionality.

Inspired by Wang et al. (2024), where this process is similar to the destruct and rebuild process of a LEGO toy, we named the fine-grained translation technique as **LEGO translation** and our system built upon it as **LEGO-Compiler**. As depicted in (b) in Figure 1, LEGO translation first breaks down large programs into manageable, self-contained blocks, analogous to LEGO pieces (**Part split**). Then these blocks are independently translated (**Part translation**) and finally recombined, enabling scalable and accurate translation of complex programs (**Part rebuild**). All these methods rely on an inherently nature in programming languages, the composability in control block level, which reflects the linearization process in compiler design (Wirth et al., 1996), where tree-structured control flow can be linearized, and therefore, composable. We have formally proved the widely applicable composability of programming languages using a constructive approach in Appendix B.

Algorithm 1 LLM-driven **Part Split** Algorithm based on Control Blocks

```

procedure SELECTCONTROLBLOCKS(function)
  blocks  $\leftarrow$   $\emptyset$ 
  deque.push_back(function)
  while deque is not empty do
    block  $\leftarrow$  deque.pop_front()
    decision  $\leftarrow$  LLMDecideSplit(block)
    if decision is "keep" then
      blocks.append(block)
    else
      subBlocks  $\leftarrow$  SplitByOutermostControl(block)
      for subBlock in subBlocks in reverse order do
        deque.push_front(subBlock)
      end for
    end if
  end while
  return blocks
end procedure

```

2.3 LEGO-COMPILER: THE FUNCTIONAL NEURAL COMPILER

We apply the LEGO translation method to the compilation domain and, guided by compilation-specific expert knowledge, design the **LEGO-Compiler**. This LLM-driven system accepts C programs as input and generates assembly code for x86, ARM, or RISC-V architectures. An overview of the LEGO-Compiler is depicted in Figure 2, primarily including the following Chain-of-Thoughts, where their detailed prompts can be found in Figure 11.

2.3.1 CONTROL FLOW ANNOTATION

The first annotation process addresses the positioning issues. Control flow in high-level languages consists of structures like if, while, for, and switch statements, which are linearized into branch-label constructs in assembly. In modern programming languages like C, most control statements (except goto) are encapsulated, meaning their generated labels remain within their scope. This property makes C programs composable at the control block level, which is formally proved in Theorem 5.

Control flow annotation is where **Part Split** is performed. Algorithm 1, which is also inspired by the composability proof algorithm in Figure 2, describes how to use LLM to split program into reasonably sized blocks. For small programs or control statements with low nesting levels, splitting may be unnecessary. For deeply nested control statements, further splitting may be required. This process maintains composability and encapsulation, ensuring correctness when rebuilding the full translation. In extreme cases, the program can be divided into basic blocks (Definition 3) or even sequential statements (Definition 2).

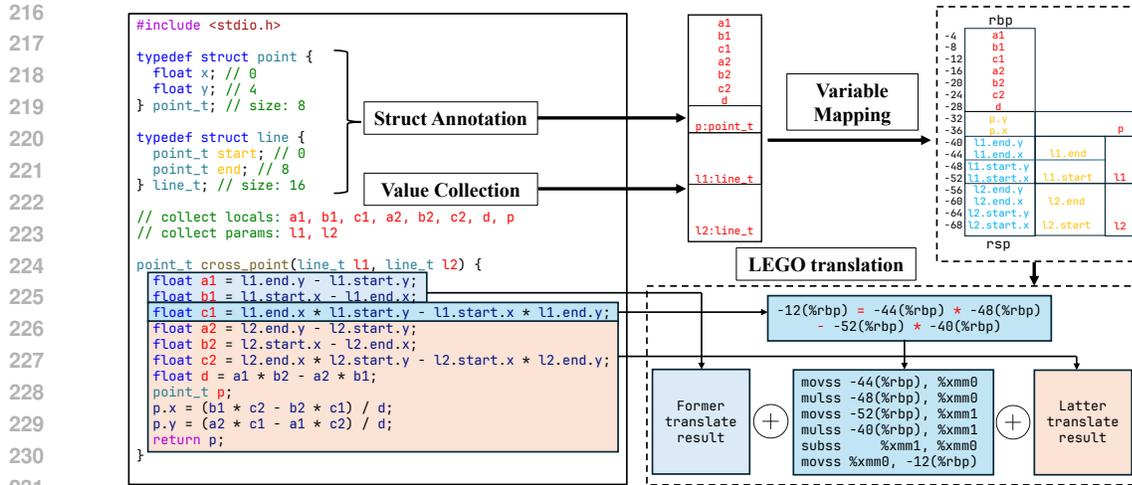


Figure 2: Workflow of **LEGO-Compiler**'s annotation-based Chain-of-Thoughts, **Struct Annotation** reasons type information from basic types, **Value Collection** finds all instances of each type, **Variable Mapping** explicitly bind variables from src language to dst language, then Part Split will split the program into composable parts, then **LEGO translate** them.

2.3.2 STRUCT ANNOTATION

The C language type system includes numerous basic types based on integer or float, their corresponding pointers, and compound types such as struct, union, and array, composed of basic types. Translating basic types and their instructions is relatively simple, as LLMs have learned this knowledge through extensive pre-training. However, for compound types, like structs (and similarly unions), the challenge arises from the infinite possible combinations of basic types.

To address this, we adopted a Chain-of-Thought approach. Instead of directly memorizing the mappings of variables from source-assembly language pairs in pretraining stage, we prompt the LLM to perform a separate thought process to reason about the memory layout of compound types based on structure, which includes size, offset, and alignment. Since compound types are ultimately composed of basic types, and LLMs understand basic type memory layouts, this pass can effectively infer the memory layout of compound types, like struct and array.

The Struct Annotation result is also verifiable using front-end tools like IntelliSense (Microsoft Corporation, 2024) or Clangd (LLVM Project, 2024b). After Struct Annotation, we obtain the symbol attributes for each type in the symbol table (LLM's context).

2.3.3 VARIABLE ANNOTATION

After type attribute inference, we need to determine where each type appears in the program and which variable identifier represents it. Our designed prompt guides the LLM to infer all declarations and arrange variable stack allocations according to their declaration order

Using x86 assembly as an example, global variables are stored in the data segment and indexed by same-name labels, providing clear binding relationships. For local variables, our method involves treating them as stored on the stack, assigning each a specific offset relative to a base address. By iterating through variable definitions, we can update these offsets relative to the base address, achieving effective stack allocation in most cases.

For compound type variables like structs, following the System V ABI (2018) for x86 assembly, access to sub-elements is achieved by adding offsets. This process is generally accurate with LLMs, relying primarily on precise binary integer arithmetic operations. After the Variable Annotation pass, we obtain a correspondence between variables in C and assembly languages, allowing for simple substitution during compilation, as illustrated in Figure 2.

270 An additional challenge is ensuring variable name uniqueness in the source program. We address
 271 this through a renaming pass at the source program level, eliminating name conflicts and ensuring
 272 variable uniqueness. This process is also verifiable through behavioral validation of the program
 273 before and after renaming.

274 2.3.4 SELF-CORRECTION THROUGH ERROR FEEDBACK

275 To address potential errors in the LLM-generated code, we implement a comprehensive self-
 276 correction mechanism. This system classifies errors into three categories: assembly semantic errors
 277 (detected by the assembler), runtime errors (identified through execution and caught by debuggers
 278 like gdb), and behavioral errors (discovered through result comparison). The error information
 279 is collected and fed back to the LLM for self-correction. Assembly semantic errors are typically
 280 straightforward to fix, while runtime errors, often caused by null pointer dereferences, are addressed
 281 by tracing instructions step-by-step to pinpoint the problematic area. Behavioral errors, being the
 282 most complex, may require multiple iterations to resolve. This iterative feedback and correction
 283 process significantly enhances the robustness and accuracy of the LLM-based compilation system,
 284 which is evaluated in the following section.

285 In general, LEGO-Compiler uses Chain-of-Thoughts in neural compilation task by either explicitly
 286 annotating the source code or generating intermediate results. These annotations and results are
 287 stored in the LLM’s context, allowing the model to integrate them effectively. Through in-context
 288 learning, LEGO-Compiler is able to perform neural compilation tasks step by step, with each step
 289 being a simpler subtask that the LLM can handle. For more details, see subsection D.1.

290 3 EXPERIMENTS

291 To evaluate the effectiveness of our LEGO-Compiler approach, we have conducted a comprehensive
 292 set of experiments using three distinct datasets: ExeBench, CoreMark, and LongFunction. Each
 293 dataset serves a specific purpose in assessing different aspects of our neural compilation method.

294 3.1 EXPERIMENTAL SETUP

295 Major parameters we have tested are listed below, All settings use a one-shot prompt to help align
 296 the format. We evaluate the **Pass@k** correctness through IO unittests, altering the following settings,
 297 note that not all combinations of experimental settings are tested due to resource constraints.

- 298 • **models:** Advanced LLMs: GPT-4o (OpenAI, 2024), Claude-3.5-sonnet (Anthropic, 2023),
 299 Deepseek-coder (Guo et al., 2024) and Mini LLMs: GPT-4o-mini, Claude-3-haiku, and
 300 Codestral-22b (AI, 2024). We also test the newest o1-preview model (OpenAI, 2024) for
 301 limited evaluation.
- 302 • **method ablation:** Direct(baseline), annotation, annotation + fixing, annotation + fixing +
 303 LEGO translation, annotation + fixing + LEGO translation + pass@k(LEGO-Compiler)
- 304 • **Temperature:** 0.0-1.0, with 0.2 step increments
- 305 • **k:** 1, 5
- 306 • **fix rounds:** 0, 1, 3
- 307 • **architecture:** x86_64, arm-v8a, riscv64, majorly on x86

308 3.2 EXEBENCH EVALUATION

309 ExeBench (Armengol-Estapé et al., 2022) is a large-scale dataset of executable C programs, each
 310 equipped with a comprehensive unittest system. We use its Real-Executable subset, initially con-
 311 taining 40k samples. After data cleaning and removing samples that couldn’t be compiled by the
 312 oracle compiler, our final test set consists of 23k samples. We utilize a test set of 500 cases randomly
 313 chosen from the full dataset for comprehensive evaluation due to resource and time constraints. To
 314 ensure the representativeness of this test set, we conduct additional evaluations on a larger subset of
 315 5,000 cases using DeepseekCoder, one of our evaluating LLMs, whose results are consistent with
 316 those obtained from the smaller subset. We evaluate ExeBench through the following methodology:

1. Translate the C program to assembly using the LLM (to generate hypothesis), where we have three methods: direct, annotation-based CoT and LEGO translation.
2. Assemble and link the hypothesis assembly to create an executable.
3. Run the executable through 10 different IO test cases provided by ExeBench.
4. Consider the translation *successful* if it passes all test cases.
5. If a translation fails, apply self-fixing with the collected error feedback, will try **fix rounds**.
6. If still unsuccessful, proceed to the next iteration in **Pass@k**, until **k** is reached.
7. Consider the translation *failed* if it doesn't pass after all configured attempts.

The overall results on ExeBench are presented in Table 1, and case-difficulty ablation results in Figure 3, where we have tested 3 advanced LLMs and 3 mini LLMs respectively with the following experimental settings. **Additionally, we compare the results with Zhang et al. (2024), the current state-of-the-art in neural compilation task, where their finetuned model is orthogonal to our work:**

- **Baseline:** Direct code translation with Pass@1 and greedy-decoding. This represents the basic neural compilation capability of LLMs based on their default pretraining results. As we can see, models vary in a large margin, advanced models outperform mini models, and Claude-series outperform GPT series, where we found GPT series are facing simple syntax failures. Additionally, DeepseekCoder and Codestral performs well in the baseline setting.
- **+Pass@k:** Altering the Temperature to 0.6 and **k** to 5. This improves greatly for trivial code syntax errors as it allows LLMs trying different styles in the assembly, however, pass@k by allowing sampling on sub-optimal choices during decoding stage can only mitigate, but not solve the inefficient pretraining of LLMs. We see large improvement for all models, where most relatively simple cases are generated correctly during this stage, however, for harder cases, and those with pretraining biases(causing the errors), Pass@k is not helpful.
- **+Feedback:** Enabling the feedback self-fixing method, which enables LLMs to self-correcting its output from assembler feedback, runtime feedback and behavioral feedback. This significantly improves those cases with pretraining biases, because by explicitly providing error feedback, LLMs will reflect on their generation and focus on solving the errors. In comparison, assembler feedback is the most efficient feedback message, because it directly points out the errors; runtime feedback is helpful as well, though LLMs need to additionally reason its actual error occurrence from the message; behavioral message, since it lacks clear information about where is wrong in the hypothesis, although it's somehow helpful, LLMs' guesses on which part is problematic are usually wrong. In general, by enabling the feedback-driven LLM self-correction method, all models get significant improvement on their Pass@5 accuracy. Typically, the advanced LLMs majorly solve the simple and medium code, while the mini LLMs will still face some problems in simple cases.
- **+CoT:** Further enabling the annotation-based CoT methods described in subsection 2.3. This helps LLMs to reason the compilation process instead of direct generation, although it requires more tokens to be consumed as it generates intermediate text and reasoning steps. As a result, a large part of hard code is successfully generated, even for mini LLMs. The possible explanation is, these mini LLMs are not sufficiently pretrained on neural compilation datasets, but are sufficient for reasoning the logic of a guided compilation process, which is given by the annotation methods. Empirical results show all models pass at least 92.2% of the testset of ExeBench, which already looks good. Additionally, the 3 advanced LLMs reach around 99% IO Accuracy, showing that except for extreme hard cases, LLMs are sufficient to translate it well with guidelines of the annotation-based Chain-of-Thoughts.
- **LEGO-Compiler(all):** Coming so far, the remaining failed cases are all difficult in at least one of their features, detailed explanations of these difficulties are in subsection D.2. By translating a managed small part of code at a time and combining these results together to form a full neural compilation, LEGO-Compiler significantly solves the difficulty due to long code length and complicated control flows, where all advanced LLMs achieve over 99% accuracy in the ExeBench testset. Since the accuracy is very impressive, we further filter the hardest 10% subset of ExeBench based on the number of basic blocks and

Table 1: ExeBench experimental results

Model	Baseline	+Pass@k	+Feedback	+CoT	LEGO-Compiler
GPT-4o	76.8%	93.4%	97.8%	99.2%	99.8%
Claude-3.5-Sonnet	92.6%	97.8%	98.6%	99.4%	100.0% ¹
DeepseekCoder ²	82.48%	87.96%	93.76%	97.36%	99.24%
GPT-4o-mini	58.8%	74.8%	86.0%	92.2%	-
Claude-3-Haiku	69.4%	81.8%	90.0%	95.8%	-
Codestral	71.8%	79.4%	88.6%	-	-
Finetuned LLM Compiler ³	91.72%	-	-	-	-

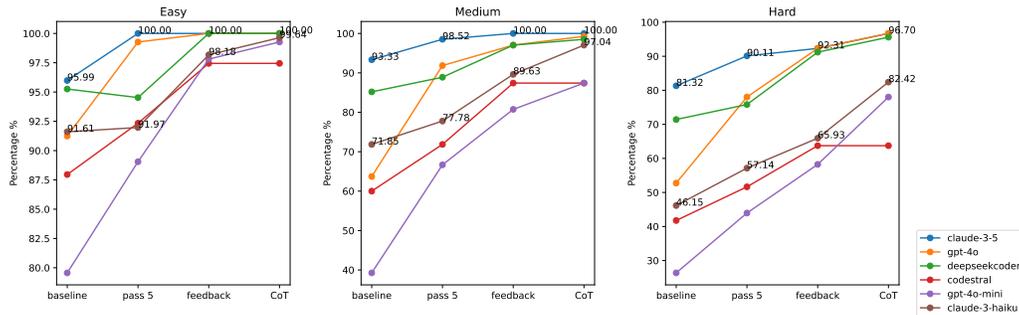


Figure 3: Ablation study: easy, medium and hard subsets ablations on ExeBench.

instructions within these blocks using the LLVM toolchain (LLVM Project, 2024a) to characterize the difficulty in ExeBench. The characterization of ExeBench and its hard subset is illustrated in Figure 6 in the appendix. Additionally, the experimental results on this hard subset, demonstrating the effectiveness of our methods, can be viewed in Table 2 also located in the appendix.

To sum up, the empirical results of our final LEGO-Compiler are suggesting the success of using LLMs for neural compilation, where the advanced LLMs solve almost all difficult cases in ExeBench(>99%), and the mini LLMs can also have over 95% accuracy.

3.3 COREMARK EVALUATION, A CASE STUDY

Previous evaluation has given promising results on neural compilation, with all methods applied, LLMs are achieving over 99% accuracy. However, considering ExeBench is a function-level compilation dataset that contains code with limited complexity, which we characterize in Figure 6, it is natural for us to think about applying our LEGO-Compiler for real-world codebases, where we choose CoreMark (Gal-On & Levy, 2012) as a case study, showing how complicated code that LLMs are capable of handling now.

CoreMark is a widely used benchmark for embedded devices, written entirely in C. It evaluates computer performance through state machine operations, linked list manipulations, and matrix computations. CoreMark consists of 40 functions, representing a complex, industry-grade codebase.

As depicted in Figure 4, the main function is one of the most complicated code in CoreMark, which contains a lot of complicated features of a C program. From another perspective, we can assert that if

¹100% test accuracy suggests LEGO-Compiler’s state-of-the-art potential, but it doesn’t ensure perfection in all scenarios, which suggests us to study harder cases as well.

²We evaluate the rest of models with a 500 subset of ExeBench, while we perform a larger scale 5000 subset of ExeBench with DeepseekCoder model, showcasing the consistency.

³We use the ExeBench results in the paper of Zhang et al. (2024) for comparison.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

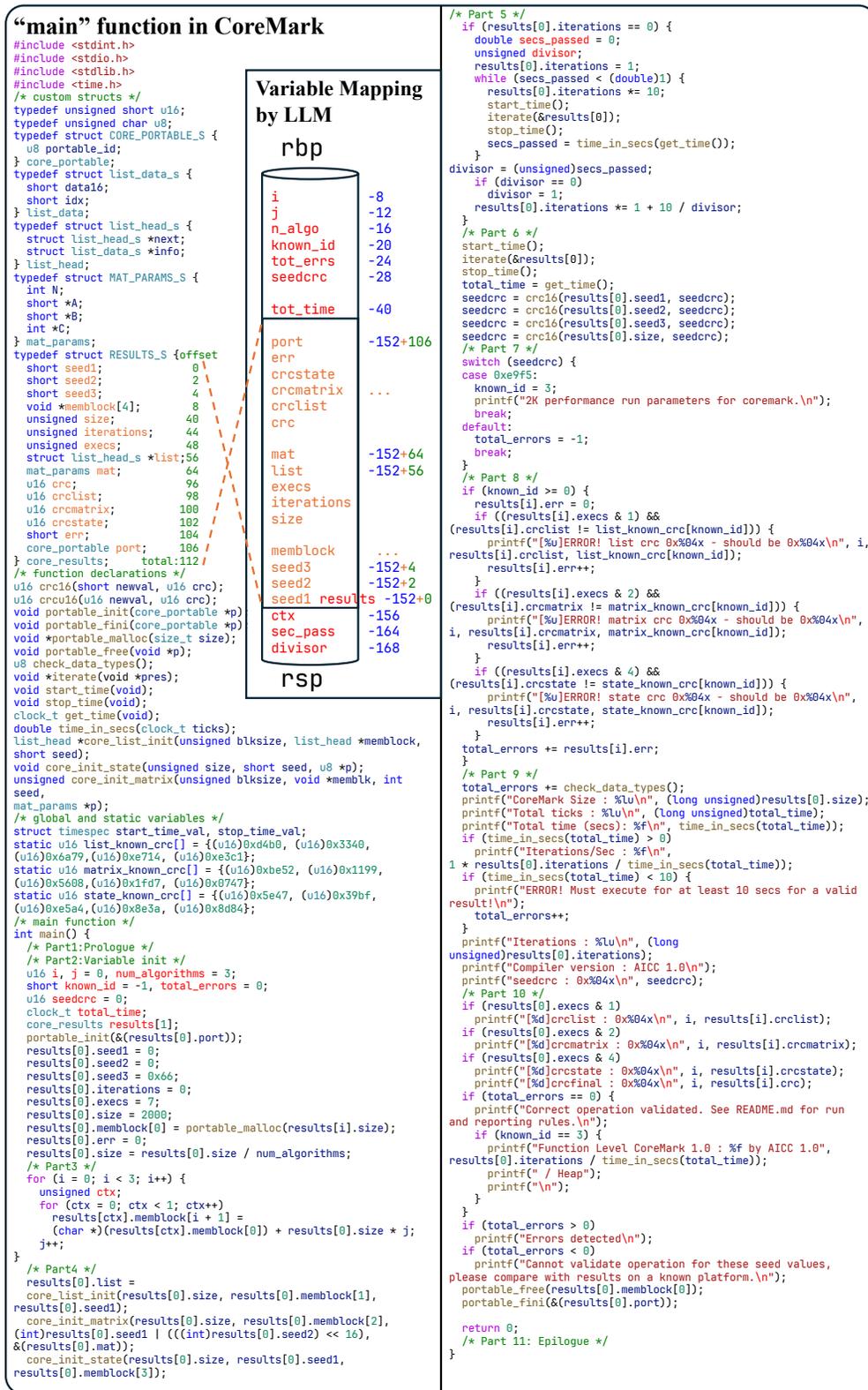


Figure 4: The CoreMark main function, one of the most difficult code we evaluated. In this figure, all CoTs are illustrated in the code annotations in color, as well as the variable mapping process.

486 CoreMark’s main function can be neural compiled, simpler code, can be neural-compiled correctly
487 in high possibility.

488
489 In general, Claude-3.5-sonnet compiles all 40 out of 40 functions correctly, both GPT-4o and
490 DeepseekCoder achieves 38, where they fail to generate the **main** function and another compli-
491 cated **core_bench_state** function. The reason for their failure is not on the complicated code control
492 structures, but on the translation of certain instructions, which can be improved with more compiler
493 knowledge taught to them. If taught with such knowledge (manually prompted), all three LLMs can
494 successfully compile the whole CoreMark, achieving functionality just as oracle compiler does.

495 496 3.4 LONGFUNCTION EVALUATION 497

498 Since the LEGO translation method significantly scales up the capability of long code translation,
499 we design LongFunction, a synthesized dataset for testing very long code translations, particularly
500 for evaluating the effectiveness of our LEGO translation method. The dataset is made up of 50 syn-
501 thesized programs, ranging from 317 to 238737 tokens in length, each program is self-contained and
502 can be compiled and run independently. We evaluate the neural compilation task on LongFunction
503 dataset for all x86_64, arm-v8a and riscv64 architectures, and the neural code translation task by
504 migrating C to Python/C++/Rust.

505 Examples of LongFunction and the evaluation details can be found in subsection D.3 due to page
506 limits. In conclusion, our proposed LEGO translation method breaks the complexity of long code,
507 boosting their capability of handling long code, where all LLMs we tested passed the whole Long-
508 Function dataset. In comparison, the current best model, o1-preview, can only maximally translate
509 a 5772 token sized case using direct translation method.

510 The results of LongFunction dataset evaluation give us a strong insight: Code, or programming
511 languages, unlike the natural languages, no matter how long they are, their complexity can be divide-
512 and-conquered into two levels: The first is on the control flow level, which combines each block of
513 code logic together to form a long and complicated code, where both current compilers and our
514 LEGO-Compiler methods can iteratively split the code into smaller and smaller, managable code
515 snippets to overcome this complexity. The second is on the handling of each code snippet, where
516 the size is not large, and the major difficulties are on how to correctly translate each statement of
517 it correctly, where we also find out that for these advanced LLMs, they perform strongly on this
518 level. Since the basic operations are also limited in programming languages, within proper size,
519 the code snippet level translation or compilation can be gradually improved to near 100% with the
520 advancement of more powerful pretrained LLMs.

521 522 4 CONCLUSION 523

524 We have presented LEGO-Compiler, a novel approach to neural compilation that leverages Large
525 Language Models (LLMs) to translate high-level programming languages into assembly code. Our
526 LEGO translation method breaks down large programs into manageable, self-contained blocks
527 through the composable nature of code, significantly extending the scalability of neural code trans-
528 lation. By incorporating a series of Chain-of-Thought stages guided by classical compiler design
529 and self-correction mechanisms, LEGO-Compiler effectively addresses key challenges in compi-
530 lation tasks, achieving significant improvements in accuracy and scalability. Our experimental re-
531 sults demonstrate the effectiveness of LEGO-Compiler, as it achieves over 99% accuracy on the
532 ExeBench dataset and fully compiles the industrial-grade CoreMark benchmark correctly. We
533 also introduce LongFunction, a new dataset designed to evaluate the translation and compilation
534 of lengthy code, demonstrating the effectiveness of the proposed LEGO translation method.

535 These findings provide important insights into the capabilities and limitations of LLMs in neural
536 compilation tasks. While our current implementation incurs higher computational costs compared
537 to traditional compilers, it offers unique advantages in adaptability and potential for rapid integration
538 of new instruction sets or language features. As LLM capabilities continue to improve, approaches
539 like LEGO-Compiler are poised to play an increasingly important role in the future of software
development and compilation, complementing and enhancing traditional compiler technologies.

REFERENCES

- 540
541
542 Mistral AI. Codestral: Hello, world! <https://mistral.ai/news/codestral/>, May
543 2024. Accessed on September 15, 2024.
- 544 Anthropic. Claude ai. <https://www.anthropic.com>, 2023. Accessed: 2024-09-14.
- 545
546 Jordi Armengol-Estapé and Michael FP O’Boyle. Learning c to x86 translation: An experiment in
547 neural compilation. *arXiv preprint arXiv:2108.07639*, 2021.
- 548 Jordi Armengol-Estapé, Jackson Woodruff, Alexander Brauckmann, José Wesley de Souza Mag-
549 alhães, and Michael FP O’Boyle. Exebench: an ml-scale dataset of executable c functions. In
550 *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pp.
551 50–59, 2022.
- 552 Jordi Armengol-Estapé, Jackson Woodruff, Chris Cummins, and Michael FP O’Boyle. Slade:
553 A portable small language model decompiler for optimized assembler. *arXiv preprint*
554 *arXiv:2305.12520*, 2023.
- 555
556 Ying Cao, Ruigang Liang, Kai Chen, and Peiwei Hu. Boosting neural networks to decompile opti-
557 mized binaries. In *Proceedings of the 38th Annual Computer Security Applications Conference*,
558 pp. 508–518, 2022.
- 559 Xinyun Chen, Chang Liu, and Dawn Song. Tree-to-tree neural networks for program translation.
560 In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*,
561 NIPS’18, pp. 2552–2562, Red Hook, NY, USA, 2018. Curran Associates Inc.
- 562 Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng,
563 Ming Liu, Bing Qin, and Ting Liu. Navigate through enigmatic labyrinth a survey of chain of
564 thought reasoning: Advances, frontiers and future, 2024. URL [https://arxiv.org/abs/
565 2309.15402](https://arxiv.org/abs/2309.15402).
- 566
567 Chris Cummins, Volker Seeker, Dejan Grubisic, Mostafa Elhoushi, Youwei Liang, Baptiste Roziere,
568 Jonas Gehring, Fabian Gloeckle, Kim Hazelwood, Gabriel Synnaeve, et al. Large language mod-
569 els for compiler optimization. *arXiv preprint arXiv:2309.07062*, 2023.
- 570 Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Roziere, Jonas Gehring, Gabriel Synnaeve,
571 and Hugh Leather. Meta large language model compiler: Foundation models of compiler opti-
572 mization, 2024. URL <https://arxiv.org/abs/2407.02524>.
- 573 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu,
574 Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning,
575 2024. URL <https://arxiv.org/abs/2301.00234>.
- 576
577 Jack J Dongarra, Piotr Luszczek, and Antoine Petit. The linpack benchmark: past, present and
578 future. *Concurrency and Computation: practice and experience*, 15(9):803–820, 2003.
- 579 Shihan Dou, Haoxiang Jia, Shenxi Wu, Huiyuan Zheng, Weikang Zhou, Muling Wu, Mingxu
580 Chai, Jessica Fan, Caishuang Huang, Yunbo Tao, Yan Liu, Enyu Zhou, Ming Zhang, Yuhao
581 Zhou, Yueming Wu, Rui Zheng, Ming Wen, Rongxiang Weng, Jingang Wang, Xunliang Cai, Tao
582 Gui, Xipeng Qiu, Qi Zhang, and Xuanjing Huang. What’s wrong with your code generated by
583 large language models? an extensive study, 2024. URL [https://arxiv.org/abs/2407.
584 06153](https://arxiv.org/abs/2407.06153).
- 585
586 Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing
587 Qin, Ting Liu, Daxin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming
588 and natural languages, 2020.
- 589 Free Software Foundation. *c++filt*. GNU Binutils, 2023. URL [https://sourceware.org/
590 binutils/docs/binutils/c_002b_002bfilt.html](https://sourceware.org/binutils/docs/binutils/c_002b_002bfilt.html). Accessed: [Insert access date
591 here].
- 592 Cheng Fu, Huili Chen, Haolan Liu, Xinyun Chen, Yuandong Tian, Farinaz Koushanfar, and Jishen
593 Zhao. Coda: An end-to-end neural program decompiler. *Advances in Neural Information Pro-
cessing Systems*, 32, 2019.

- 594 Shay Gal-On and Markus Levy. Exploring coremark a benchmark maximizing simplicity and effi-
595 cacy. *The Embedded Microprocessor Benchmark Consortium*, 2012.
596
- 597 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
598 Bi, Y. Wu, Y.K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the
599 large language model meets programming – the rise of code intelligence, 2024. URL <https://arxiv.org/abs/2401.14196>.
600
- 601 Zifan Carl Guo and William S. Moses. Enabling transformers to understand low-level programs.
602 In *2022 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1–9, 2022. doi:
603 10.1109/HPEC55821.2022.9926313.
604
- 605 John L Gustafson and Quinn O Snell. Hint: A new way to measure computer performance. In
606 *Proceedings of the Twenty-Eighth Annual Hawaii International Conference on System Sciences*,
607 volume 2, pp. 392–401. IEEE, 1995.
- 608 Ali Reza Ibrahimzada. Program decomposition and translation with static analysis. In *Proceedings*
609 *of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion*
610 *Proceedings, ICSE-Companion '24*, pp. 453–455, New York, NY, USA, 2024. Association for
611 Computing Machinery. ISBN 9798400705021. doi: 10.1145/3639478.3641226. URL <https://doi.org/10.1145/3639478.3641226>.
612
- 613 Siddharth Jha, Lutfi Eren Erdogan, Sehoon Kim, Kurt Keutzer, and Amir Gholami. Characterizing
614 prompt compression methods for long context inference. In *Workshop on Efficient Systems for*
615 *Foundation Models II @ ICML2024*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=vs6CCDuK71)
616 [vs6CCDuK71](https://openreview.net/forum?id=vs6CCDuK71).
617
- 618 Mingsheng Jiao, Tingrui Yu, Xuan Li, Guanjie Qiu, Xiaodong Gu, and Beijun Shen. On the eval-
619 uation of neural code translation: Taxonomy and benchmark, 2023. URL <https://arxiv.org/abs/2308.08961>.
620
- 621 Bonan Kou, Shengmai Chen, Zhijie Wang, Lei Ma, and Tianyi Zhang. Do large language models
622 pay similar attention like human programmers when generating code? *Proc. ACM Softw. Eng.*, 1
623 (FSE), jul 2024. doi: 10.1145/3660807. URL <https://doi.org/10.1145/3660807>.
624
- 625 Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev.
626 In search of needles in a 11m haystack: Recurrent memory finds what llms miss, 2024. URL
627 <https://arxiv.org/abs/2402.10790>.
- 628 Mosh Levy, Alon Jacoby, and Yoav Goldberg. Same task, more tokens: the impact of input length
629 on the reasoning performance of large language models. In Lun-Wei Ku, Andre Martins, and
630 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computa-*
631 *tional Linguistics (Volume 1: Long Papers)*, pp. 15339–15353, Bangkok, Thailand, August 2024.
632 Association for Computational Linguistics. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.acl-long.818)
633 [acl-long.818](https://aclanthology.org/2024.acl-long.818).
- 634 Zhuowan Li, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Retrieval aug-
635 mented generation or long-context llms? a comprehensive study and hybrid approach, 2024.
636 URL <https://arxiv.org/abs/2407.16833>.
637
- 638 J. Liu, F. Zhang, X. Zhang, Z. Yu, L. Wang, Y. Zhang, and B. Guo. hmcodetrans: Human-machine
639 interactive code translation. *IEEE Transactions on Software Engineering*, 50(05):1163–1181,
640 may 2024a. ISSN 1939-3520. doi: 10.1109/TSE.2024.3379583.
- 641 Yue Liu, Chakkrit Tantithamthavorn, Yonghui Liu, and Li Li. On the reliability and explainability of
642 language models for program generation. *ACM Trans. Softw. Eng. Methodol.*, 33(5), jun 2024b.
643 ISSN 1049-331X. doi: 10.1145/3641540. URL <https://doi.org/10.1145/3641540>.
- 644 LLVM Project. *The LLVM Compiler Infrastructure*. LLVM Foundation, 2024a. URL <https://llvm.org>. Version 18.1.8.
645
646
- 647 LLVM Project. Clangd: C/c++ language server. <https://clangd.llvm.org/>, 2024b. Ac-
cessed: 2024-09-14.

- 648 Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin
649 Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. Codexglue: A machine learning benchmark
650 dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- 651
652 Microsoft Corporation. Intellisense in visual studio code. <https://code.visualstudio.com/docs/editor/intellisense>, 2024. Accessed: 2024-09-14.
- 653
654 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
655 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In
656 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*
657 *on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United
658 Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/
659 2022.emnlp-main.759. URL <https://aclanthology.org/2022.emnlp-main.759>.
- 660
661 Christian Munley, Aaron Jarmusch, and Sunita Chandrasekaran. Llm4vv: Developing llm-driven
662 testsuite for compiler validation. *Future Generation Computer Systems*, 2024.
- 663
664 nfnit. Ansibench: A selection of ansi c benchmarks and programs useful as benchmarks, 2024.
665 URL <https://github.com/nfnit/ansibench>. Accessed: 2024-11-22.
- 666
667 Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N. Nguyen. Divide-and-conquer approach for
668 multi-phase statistical migration for source code. In *Proceedings of the 30th IEEE/ACM Inter-*
669 *national Conference on Automated Software Engineering, ASE '15*, pp. 585–596. IEEE Press,
670 2015a. ISBN 9781509000241. doi: 10.1109/ASE.2015.74. URL <https://doi.org/10.1109/ASE.2015.74>.
- 671
672 Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N. Nguyen. Divide-and-conquer approach
673 for multi-phase statistical migration for source code (t). In *2015 30th IEEE/ACM Inter-*
674 *national Conference on Automated Software Engineering (ASE)*, pp. 585–596, 2015b. doi:
10.1109/ASE.2015.74.
- 675
676 Theo X. Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-
677 Lezama. Is self-repair a silver bullet for code generation? In *The Twelfth International Confer-*
678 *ence on Learning Representations*, 2024. URL <https://openreview.net/forum?id=y0GJXRungR>.
- 679
680 OpenAI. Chatgpt: Optimizing language models for dialogue. *OpenAI*, 2023. URL <https://openai.com/research/chatgpt>.
- 681
682 OpenAI. Gpt-4o system card. <https://openai.com/index/gpt-4o-system-card/>,
683 August 2024. Accessed on September 15, 2024.
- 684
685 OpenAI. Openai o1 system card. Technical report, OpenAI, September 2024. URL <https://cdn.openai.com/o1-system-card.pdf>. Accessed on September 15, 2024.
- 686
687 OpenAI. Video generation models as world simulators. [https://openai.com/index/
688 video-generation-models-as-world-simulators/](https://openai.com/index/video-generation-models-as-world-simulators/), 2024. Technical Report.
- 689
690 OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
691 cia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red
692 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavar-
693 ian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner,
694 Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim
695 Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany
696 Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek
697 Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
698 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
699 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling,
700 Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko
701 Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Chris-
tian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon,
Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei

- 702 Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke,
703 Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu,
704 Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang,
705 Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaf-
706 tan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan,
707 Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie
708 Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstan-
709 tinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan
710 Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz
711 Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning,
712 Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew,
713 Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok
714 Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan
715 Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Rei-
716 chiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang,
717 Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista
718 Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam
719 Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
720 Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power,
721 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Fran-
722 cis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli,
723 Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman,
724 Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam,
725 Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Ben-
726 jamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya
727 Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Eliza-
728 beth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone,
729 Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang,
730 Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian
731 Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren
732 Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
733 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
734 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2023.
- 735
736 Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang.
737 Automatically correcting large language models: Surveying the landscape of diverse automated
738 correction strategies. *Transactions of the Association for Computational Linguistics*, 12:484–506,
739 2024. doi: 10.1162/tacl.a.00660. URL [https://aclanthology.org/2024.tacl-1.](https://aclanthology.org/2024.tacl-1.27)
740 27.
- 741
742 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
743 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
744 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 745
746 Baptiste Rozière, Marie-Anne Lachaux, Lowik Chanussot, and Guillaume Lample. Unsupervised
747 translation of programming languages. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Had-
748 sell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Process-*
749 *ing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS*
750 *2020, December 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2020/hash/ed23fbf18c2cd35f8c7f8de44f85c08d-Abstract.html)
751 [paper/2020/hash/ed23fbf18c2cd35f8c7f8de44f85c08d-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/ed23fbf18c2cd35f8c7f8de44f85c08d-Abstract.html).
- 752
753 Baptiste Roziere, Jie M Zhang, Francois Charton, Mark Harman, Gabriel Synnaeve, and Guil-
754 laume Lample. Leveraging automated unit tests for unsupervised code translation. *arXiv preprint*
755 *arXiv:2110.06773*, 2021.
- 756
757 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
758 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code.
759 *arXiv preprint arXiv:2308.12950*, 2022.

- 756 Xinyu She, Yue Liu, Yanjie Zhao, Yiling He, Li Li, Chakkrit Tantithamthavorn, Zhan Qin, and
757 Haoyu Wang. Pitfalls in language models for code intelligence: A taxonomy and survey. *arXiv*
758 *preprint arXiv:2310.17903*, 2023.
- 759
- 760 Mingyang Song, Mao Zheng, and Xuan Luo. Can many-shot in-context learning help long-context
761 llm judges? see more, judge better!, 2024. URL <https://arxiv.org/abs/2406.11629>.
- 762
- 763 System V ABI. System v application binary interface: Amd64 architecture processor supple-
764 ment (with lp64 and ilp32 programming models) version 1.0. Technical report, The Santa Cruz
765 Operation, Inc., 2018. URL [https://github.com/hjl-tools/x86-psABI/wiki/
766 x86-64-psABI-1.0.pdf](https://github.com/hjl-tools/x86-psABI/wiki/x86-64-psABI-1.0.pdf).
- 767
- 768 Marc Szafraniec, Baptiste Rozière, Hugh Leather, Patrick Labatut, François Charton, and Gabriel
769 Synnaeve. Code translation with compiler representations. In *The Eleventh International Confer-
770 ence on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
771 2023. URL <https://openreview.net/pdf?id=XomEU3eNeSQ>.
- 772
- 773 Bowen Tan, Zichao Yang, Maruan Al-Shedivat, Eric Xing, and Zhiting Hu. Progressive gener-
774 ation of long text with pretrained language models. In Kristina Toutanova, Anna Rumshisky,
775 Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy
776 Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North Amer-
777 ican Chapter of the Association for Computational Linguistics: Human Language Technologies*,
778 pp. 4313–4324, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/
2021.naacl-main.341. URL <https://aclanthology.org/2021.naacl-main.341>.
- 779
- 780 Kiran Vodrahalli, Santiago Ontanon, Nilesh Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna,
781 Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, Rohan Anil, Ethan Dyer, Siamak
782 Shakeri, Roopali Vij, Harsh Mehta, Vinay Ramasesh, Quoc Le, Ed Chi, Yifeng Lu, Orhan Firat,
783 Angeliki Lazaridou, Jean-Baptiste Lespiau, Nithya Attaluri, and Kate Olszewska. Michelangelo:
784 Long context evaluations beyond haystacks via latent structure queries, 2024. URL <https://arxiv.org/abs/2409.12640>.
- 785
- 786 Boshi Wang, Xiang Deng, and Huan Sun. Iteratively prompt pre-trained language models for
787 chain of thought. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings
788 of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2714–2730,
789 Abu Dhabi, United Arab Emirates, December 2022a. Association for Computational Linguis-
790 tics. doi: 10.18653/v1/2022.emnlp-main.174. URL [https://aclanthology.org/2022.
791 emnlp-main.174](https://aclanthology.org/2022.emnlp-main.174).
- 792
- 793 Haiming Wang, Huajian Xin, Chuanyang Zheng, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing
794 Xiong, Han Shi, Enze Xie, Jian Yin, Zhenguo Li, and Xiaodan Liang. LEGO-prover: Neural
795 theorem proving with growing libraries. In *The Twelfth International Conference on Learning
796 Representations*, 2024. URL <https://openreview.net/forum?id=3f5PALef5B>.
- 797
- 798 Xin Wang, Yasheng Wang, Yao Wan, Fei Mi, Yitong Li, Pingyi Zhou, Jin Liu, Hao Wu, Xin Jiang,
799 and Qun Liu. Compilable neural code generation with compiler feedback. In Smaranda Muresan,
800 Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Computational
801 Linguistics: ACL 2022*, pp. 9–19, Dublin, Ireland, May 2022b. Association for Computational
802 Linguistics. doi: 10.18653/v1/2022.findings-acl.2. URL [https://aclanthology.org/
803 2022.findings-acl.2](https://aclanthology.org/2022.findings-acl.2).
- 804
- 805 Yue Wang, Weishi Wang, Shafiq Joty, and Steven C. H. Hoi. Codet5: Identifier-aware unified pre-
806 trained encoder-decoder models for code understanding and generation, 2021.
- 807
- 808 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V
809 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models.
In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in
Neural Information Processing Systems*, volume 35, pp. 24824–24837. Curran Associates, Inc.,
2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/
file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf).

- 810 Yuanbo Wen, Qi Guo, Qiang Fu, Xiaqing Li, Jianxing Xu, Yanlin Tang, Yongwei Zhao, Xing Hu,
811 Zidong Du, Ling Li, et al. Babeltower: Learning to auto-parallelized program translation. In
812 *International Conference on Machine Learning*, pp. 23685–23700. PMLR, 2022.
- 813
814 Niklaus Wirth, Niklaus Wirth, Niklaus Wirth, Suisse Informaticien, and Niklaus Wirth. *Compiler*
815 *construction*, volume 1. Addison-Wesley Reading, 1996.
- 816 Yiqing Xie, Atharva Naik, Daniel Fried, and Carolyn Rose. Data augmentation for code trans-
817 lation with comparable corpora and multiple references. In Houda Bouamor, Juan Pino,
818 and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP*
819 *2023*, pp. 13725–13739, Singapore, December 2023. Association for Computational Linguistics.
820 doi: 10.18653/v1/2023.findings-emnlp.917. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-emnlp.917)
821 [findings-emnlp.917](https://aclanthology.org/2023.findings-emnlp.917).
- 822 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in c com-
823 pilers. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language*
824 *Design and Implementation*, PLDI ’11, pp. 283–294, New York, NY, USA, 2011. Associa-
825 tion for Computing Machinery. ISBN 9781450306638. doi: 10.1145/1993498.1993532. URL
826 <https://doi.org/10.1145/1993498.1993532>.
- 827
828 Zhen Yang, Fang Liu, Zhongxing Yu, Jacky Wai Keung, Jia Li, Shuo Liu, Yifan Hong, Xiaoxue Ma,
829 Zhi Jin, and Ge Li. Exploring and unleashing the power of large language models in automated
830 code translation. *Proc. ACM Softw. Eng.*, 1(FSE), jul 2024. doi: 10.1145/3660778. URL <https://doi.org/10.1145/3660778>.
- 831
832 Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu, Bingchao Wu, Bei Guan, Wang Yongji, and
833 Jian-Guang Lou. Large language models meet NL2Code: A survey. In Anna Rogers, Jordan
834 Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Associ-*
835 *ation for Computational Linguistics (Volume 1: Long Papers)*, pp. 7443–7464, Toronto, Canada,
836 July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.411. URL
837 <https://aclanthology.org/2023.acl-long.411>.
- 838 Shuoming Zhang, Jiacheng Zhao, Chunwei Xia, Zheng Wang, Yunji Chen, and Huimin Cui. In-
839 troducing compiler semantics into large language models as programming language transla-
840 tors: A case study of C to x86 assembly. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung
841 Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 996–
842 1011, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL
843 <https://aclanthology.org/2024.findings-emnlp.55>.
- 844 Albert Ziegler, Eirini Kalliamvakou, X. Alice Li, Andrew Rice, Devon Rifkin, Shawn Simister,
845 Ganesh Sittampalam, and Edward Aftandilian. Measuring github copilot’s impact on productivity.
846 *Commun. ACM*, 67(3):54–63, February 2024. ISSN 0001-0782. doi: 10.1145/3633453. URL
847 <https://doi.org/10.1145/3633453>.
- 848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A RELATED WORK

A.1 CODE TRANSLATION

Code Translation has evolved from traditional statistical methods (Nguyen et al., 2015b) to neural-based approaches that capture programming language structures (Chen et al., 2018). Current neural code translation researches can be majorly categorized to two types: learning-based translators (Rozière et al., 2020; Roziere et al., 2021; Wen et al., 2022) and pre-trained language models (Feng et al., 2020; Wang et al., 2021; Lu et al., 2021; Rozière et al., 2022; OpenAI et al., 2023; Anthropic, 2023). The former majorly studies the scarcity of parallel corpora (Xie et al., 2023) and develops unsupervised learning methods to overcome it. The latter using Large Language Models’ vast pretrained knowledge, can also perform code translations well without training (Yang et al., 2024; Liu et al., 2024a).

Analysis of neural code translation is equally crucial. Studies have examined common pitfalls in language models for code intelligence (She et al., 2023; Jiao et al., 2023), investigated the reliability and explainability of these models in automated program generation (Liu et al., 2024b), and the attention paid by LLM during code generation that differs from human (Kou et al., 2024).

As for **compilation related translations**, Armengol-Estapé & O’Boyle (2021) first gives a try on neural compilation. Guo & Moses (2022) further studies on C-to-LLVM IR translation. However, they only perform limited investigations on the methods, and their results are still preliminary. There are also works on the reverse decompilation process (Fu et al., 2019; Cao et al., 2022; Armengol-Estapé et al., 2023) and works on code optimizations (Cummins et al., 2023; 2024). *The most related work is Zhang et al. (2024), which achieves state-of-the-art 91% Pass@1 accuracy compared to non-LLM methods on the C-x86 task using a finetuned CodeLlama model, which our work surpasses. Besides, their approach relies on compiler-generated bilingual corpora, while our methods can effectively eliminate such dependency by reasoning the steps of how a compiler works.*

Finally, the **breakdown of neural code translation** is also less studied, Nguyen et al. (2015a) first breaks the translation of Java-C# into syntaxemes level to lower the translation difficulty in SMTs. Our work uses similar divide-and-conquer methodology to breakdown a large long code into manageable control block parts, then LLMs can translate these parts separately with the aid of necessary context and combine their results into a large, complete and coherent translation.

A.2 OTHER RELATED WORK

LLM self-repair. Recent research has focused extensively on enhancing LLMs’ self-correction capabilities. Several studies closely related to our work deserve mention. A comprehensive survey by Pan et al. (2024) thoroughly examined methods for leveraging feedback to autonomously improve LLM outputs. Wang et al. (2022b) first uses compiler feedback for better code generation, and Dou et al. (2024) establishes the syntax-runtime-functional bug type taxonomy and builds corresponding self-repair pipelines for code. Our work is their natural extensions to neural compilation scenario. While Olausson et al. (2024) investigated the limitations of self-repair mechanisms in code generation, our findings diverge significantly. Contrary to their conclusions, we discovered that self-repair serves as a highly effective solution in the neural compilation process, particularly when incorporating syntax feedback and runtime feedback.

In-context learning and Chain-of-Thoughts. LLMs are able to in-context learn via inference alone by providing few shots of demonstration then predicting on new inputs (Min et al., 2022; Dong et al., 2024). Thus customized Chain-of-Thoughts (Wei et al., 2022; Chu et al., 2024) can guide LLMs to perform complicated reasoning (Wang et al., 2022a; Song et al., 2024), which is the cornerstone of our work. More specifically, Levy et al. (2024) reveals the degradation of LLMs’ performance for long context, and validate the effectiveness on using Chain-of-Thoughts to mitigate. We found similar results in code translation/compilation tasks. However, our proposed **LEGO translations** method can significantly mitigate such degradation as it turns a long context direct translation into multiple composable, shorter ones that LLMs can handle.

Generation Scalability and Long Context Learning. Except for code translation, many LLM-based method will fall into scalability problems since larger inputs are not well trained like the smaller ones. So methods to extend LLMs scalability remain an interesting study. For example,

in order to coherently generate long passages of text, Tan et al. (2021) proposes a multi-staged keyword-first progressive method to improve it significantly, where our work shares a similar insight. Li et al. (2024) introduces a self-route method to dynamically choose the usage of RAG or fully in-context, balancing the cost and performance in long-context scenario, which inspires us to use an analyze-first, then-CoT approach.

Needle-in-the-haystack experiment (Kuratov et al., 2024) is a well-known test for testing LLMs’ capability for long context, however, it only requires simple reasoning on the needle part, where the test is not complicated enough. There are more works evaluating the long-context learning capabilities of LLMs. Vodrahalli et al. (2024) examines LLMs with their proposed Latent Structure Queries evaluation, which aims to chisel away irrelevant information in the context, revealing a latent structure in the context, which provide a stronger signal of long-context language model capabilities. Prompt compression is another useful method to improve the long-context inference capabilities (Jha et al., 2024), which is widely used for retrieval-augmented generation(RAG) systems by compressing the long contexts. Our work and its broader area: neural compilation/translation in large codebase, could serve as another useful real-world application for long-context inference.

B COMPOSABILITY OF C-LIKE LANGUAGE CONSTRUCTS

B.1 DEFINITIONS AND LANGUAGE STRUCTURE

We define a simplified C-like language structure using the following EBNF-inspired grammar:

```

block: '{' (blockItem) * '}';
blockItem: decl | stmt;
stmt:
    lVal '=' exp ';'           # assignStmt
    | exp ';'                 # exprStmt
    | 'goto' label ';'        # gotoStmt
    | ';'                     # blankStmt
    | block                   # blockStmt
    | IF '(' exp ')' stmt (ELSE stmt)? # ifStmt
    | WHILE '(' exp ')' stmt  # whileStmt
    | FOR '(' stmt exp ';' stmt ')' stmt # forStmt
    | SWITCH '(' stmt ')' stmt # switchStmt
    | BREAK ';'              # breakStmt
    | CONTINUE ';'          # continueStmt
    | RETURN (exp)? ';'     # returnStmt;

```

We derived from the grammar that describes C-like language to form the following definitions. Also for simplicity purposes, we omit the slight differences between **decl**, **stmt** and **exp**.

Definition 2 (Basic Statement). *A basic statement is a statement that does not contain any other statements within its structure. This includes assignStmt, exprStmt, gotoStmt, blankStmt, breakStmt, continueStmt, and returnStmt. We first exclude gotoStmt for the main proof for simplicity.*

Definition 3 (Basic Block). *A basic block is a sequence of consecutive basic statements as defined in Definition 2, in which flow of control enters at the beginning and leaves at the end without halt or possibility of branching except at the end.*

Definition 4 (Control Block). *A control block is a code snippet that reflects a complete control structure, such as for(;){}, if(){else{}}, while(){}, do{}while(), or switch(){case:...}. Each subpart of a control block can be other control blocks or basic blocks as defined in Definition 3.*

Definition 5. *A basic control block is an innermost control block (Definition 4) where each of its subparts contains only basic blocks as defined in Definition 3.*

Definition 6 (Compound Control Block). *A compound control block is a control block (Definition 4) that contains at least one subpart that is not a basic block (Definition 3), but rather another control block as defined in Definition 4.*

Definition 7 (Translation Function and Valid Translations). *Let \mathcal{T} be the set of all valid translation functions from SRC to DST, where SRC is the source language (our C-like language) and DST is the destination language (e.g., x86 assembly).*

972 Formally, $\mathcal{T} = \{T \mid T : SRC \rightarrow DST\}$ such that for any $T \in \mathcal{T}$ and any $stmt \in SRC$:

973
974 1. $T(stmt) \in DST$ 2. $T(stmt)$ preserves the semantics of $stmt$

975 A translation function $T \in \mathcal{T}$ maps each construct in the source language to one or more constructs
976 in the destination language while preserving the program's behavior.

977 **Definition 8** (Translation Composability). Let (SRC, \circ) be the source language with concatenation
978 operation \circ , and (DST, \cdot) be the destination language with concatenation operation \cdot . Let \mathcal{T} be the
979 set of valid translation functions as defined in Definition 7.

980 Translation composability holds if and only if:

$$981 \quad \exists T \in \mathcal{T} : \forall P_1, P_2 \in SRC, T(P_1 \circ P_2) \equiv T(P_1) \cdot T(P_2)$$

982
983 Where:

- 984 • $T : SRC \rightarrow DST$ is a translation function
- 985 • \equiv denotes semantic equivalence, preserving both control flow and data flow
- 986 • $\circ : SRC \times SRC \rightarrow SRC$ is the concatenation operation in the source language
- 987 • $\cdot : DST \times DST \rightarrow DST$ is the concatenation operation in the destination language

993 B.2 COMPOSABILITY OF BASIC STATEMENTS

994
995 **Theorem 1** (Composability of Basic Statements). For any two basic statements $stmt_1$ and $stmt_2$
996 in SRC , as defined in Definition 2, their translation is composable: $T(stmt_1 \circ stmt_2) \equiv T(stmt_1) \cdot$
997 $T(stmt_2)$

998
999 *Proof.* We prove this for all combinations of assignment statements and expression statements. The
1000 proof considers control flow preservation, data flow preservation, and independence of translation.
1001 Other basic statements (blank, return, etc.) trivially maintain composability as they do not affect
1002 control or data flow when composed with other basic statements. \square

1003 **Example B.1.** This example illustrates the composability of basic statements as defined in Defini-
1004 tion 2 and proved in Theorem 1.

1005 Consider the following sequence of basic statements:

```
1007 a = b + 3; // stmt_1
1008 b = a - 1; // stmt_2
```

1009 The translation of these statements might look like:

```
1010 T(stmt_1):
1011   mov eax, [b]
1012   add eax, 3
1013   mov [a], eax
```

```
1014 T(stmt_2):
1015   mov eax, [a]
1016   sub eax, 1
1017   mov [b], eax
```

1018
1019 These translations are composable because:

1020
1021 1. Control Flow: The order of execution is preserved ($stmt_1$ then $stmt_2$). 2. Data Flow: The value
1022 of 'a' computed in $stmt_1$ is correctly used in $stmt_2$. 3. Independence: The translation of $stmt_2$
1023 does not depend on how $stmt_1$ was translated, only on its effect (the value of 'a').
1024

1025 Therefore, $T(stmt_1 \circ stmt_2) \equiv T(stmt_1) \cdot T(stmt_2)$, demonstrating composability.

1026 Example B.1 illustrates that even when statements have data dependencies, their translations remain
 1027 composable as long as the order of operations is preserved. Similar proof of composability can be
 1028 made for all stmts within a basic block (Definition 3).

1029 1030 B.3 COMPOSABILITY OF BASIC CONTROL STRUCTURES

1031
1032 **Theorem 2** (Composability of Basic Control Structures). *Basic control structures (if-else, for, while,*
 1033 *do-while, switch-case), where all their components are basic blocks as defined in Definition 3, are*
 1034 *composable under the translation function T as defined in Definition 7.*

1035
1036 *Proof.* We will prove this for each basic control structure:

1037 1. For Loop:

1038 Let B_{init} , B_{cond} , B_{incr} , and B_{body} be the basic blocks for init, cond, incr, and body respectively.

1039 Translation structure:

```
1040 T(basic_for_loop) :
1041     T(B_init)
1042 loop_start :
1043     T(B_cond)
1044     jz loop_end
1045     T(B_body)
1046     T(B_incr)
1047     jmp loop_start
1048 loop_end :
```

1049
1050
1051 1. Control Flow Preservation: The structure of jump instructions preserves the original control
 1052 flow. 2. Data Flow Preservation: The order of operations within and between blocks is maintained.
 1053 3. Composability: $T(\text{basic_for_loop}) \equiv T(B_{init}) \cdot T(B_{cond}) \cdot T(B_{body}) \cdot T(B_{incr})$, where \cdot
 1054 represents concatenation with appropriate jump instructions.

1055 Therefore, the basic for loop is composable under T . Similar proofs can be constructed for other
 1056 basic control structures. \square

1057
1058 2. If-Else Statement: Let B_{cond} , B_{then} , and B_{else} be the basic blocks for condition, then-branch,
 1059 and else-branch respectively.

1060 Translation structure:

```
1061 T(basic_if_else) :
1062     T(B_cond)
1063     jz else_label
1064     T(B_then)
1065     jmp end_label
1066 else_label :
1067     T(B_else)
1068 end_label :
```

1070 Control flow and data flow preservation follow similarly to the for loop case.

1071
1072 3. While Loop: Let B_{cond} and B_{body} be the basic blocks for condition and body respectively.

1073 Translation structure:

```
1074 T(basic_while) :
1075 loop_start :
1076     T(B_cond)
1077     jz loop_end
1078     T(B_body)
1079     jmp loop_start
```

1080 loop_end:
1081
1082 4. Do-While Loop: Let B_{body} and B_{cond} be the basic blocks for body and condition respectively.
1083 Translation structure:
1084
1085 T(basic_do_while):
1086 loop_start:
1087 T(B_{body})
1088 T(B_{cond})
1089 jnz loop_start
1090
1091 5. Switch-Case Statement: Let B_{expr} be the basic block for the switch expression, and
1092 B_1, B_2, \dots, B_n be the basic blocks for each case.
1093 Translation structure:
1094
1095 T(basic_switch):
1096 T(B_{expr})
1097 cmp result, case1_value
1098 je case1_label
1099 cmp result, case2_value
1100 je case2_label
1101 ...
1102 jmp default_label
1103 case1_label:
1104 T(B_1)
1105 // No break implies fall-through
1106 case2_label:
1107 T(B_2)
1108 ...
1109 default_label:
1110 T(B_n)
1111 end_switch:
1112
1113 For all these structures, control flow is preserved by the appropriate use of jump instructions, and
1114 data flow is maintained by the sequential execution of basic blocks. The translation of each structure
1115 is a composition of its basic block translations, proving composability.
1116
1117 **Theorem 3** (Composability of Break and Continue Statements). *Break and continue statements,*
1118 *which are basic statements as per Definition 2, are composable within their respective control struc-*
1119 *tures when proper loop depth tracking is maintained.*
1120
1121 *Proof.* Let $loop_depth$ be a counter maintained during translation to track nested loop levels.
1122
1123 1. Break Statement: Translation structure:
1124
1125 T(break):
1126 jmp loop_end_label_depth
1127
1128 Where $loop_end_label_depth$ corresponds to the end of the current loop at depth $loop_depth$.
1129
1130 2. Continue Statement: Translation structure:
1131
1132 T(continue):
1133 jmp loop_continue_label_depth
1134
1135 Where $loop_continue_label_depth$ corresponds to the continuation point of the current loop at depth
1136 $loop_depth$.
1137
1138 Control flow is preserved by jumping to the appropriate label based on the current loop depth. Data
1139 flow is trivially preserved as these statements do not modify data.

```

1134
1135 Algorithm 2 Iterative Bottom-Up Composability Proof Algorithm
1136
1137 procedure PROVECOMPOSABILITY(Program  $P$ )
1138    $blocks \leftarrow$  DecomposeIntoOutermostControlBlocks( $P$ )           ▷ Initial decomposition
1139    $to\_process \leftarrow$  new Deque()
1140   for each  $block$  in  $blocks$  do
1141      $to\_process.PushBack(block)$                                      ▷ Initialize processing queue
1142   end for
1143   while  $to\_process$  is not empty do
1144      $current\_block \leftarrow to\_process.PopFront()$                  ▷ Handle first unhandled block
1145     if IsBasicBlock( $current\_block$ ) then
1146       continue                                                   ▷ Do nothing
1147     else if IsControlStructure( $current\_block$ ) then
1148        $sub\_blocks \leftarrow$  SplitControlStructure( $current\_block$ )
1149       for each  $sub\_block$  in  $sub\_blocks$  in reverse order do
1150          $to\_process.PushFront(sub\_block)$                          ▷ Handle sub-blocks in original order
1151       end for
1152     else
1153       return  $P$  is not composable                                 ▷ Unrecognized structure
1154     end if
1155   end while
1156   return  $P$  is composable
1157 end procedure
1158 function SPLITCONTROLSTRUCTURE(Block  $b$ )
1159   if  $b$  is a For Loop then
1160     return SplitForLoop( $b$ )
1161   else if  $b$  is an If-Else structure then
1162     return SplitIfElse( $b$ )
1163   else
1164     return SplitOtherControlStructure( $b$ )                       ▷ Extensible for other structures
1165   end if
1166 end function
1167 function SPLITFORLOOP(ForLoop  $f$ )
1168   return [  $f.init$ ,  $f.ForBodyLabel$ ,  $f.cmp$ , ConditionalJump( $f.ForEndLabel$ ),
1169    $f.body$ ,  $f.incr$ , UnconditionalJump( $f.ForBodyLabel$ ),  $f.ForEndLabel$  ]
1170 end function
1171 function SPLITIFELSE(IfElse  $i$ )
1172   return [  $i.cmp$ , ConditionalJump( $i.ElseLabel$ ),  $i.then\_body$ ,
1173   UnconditionalJump( $i.EndIfLabel$ ),  $i.ElseLabel$ ,  $i.else\_body$ ,  $i.EndIfLabel$  ]
1174 end function

```

1172
1173
1174
1175
1176
1177

The composability of these statements within their containing loops is maintained because: a) They generate a single jump instruction that integrates with the loop’s control flow. b) The loop depth tracking ensures the jump targets the correct loop level in nested structures. \square

1178
1179

B.4 COMPOSABILITY OF COMPLEX STRUCTURES

1180

Definition 9 (Composable Control Block). *A composable control block is either:*

1181
1182
1183
1184
1185
1186
1187

- A basic block as defined in Definition 3, or
- A basic control structure as proved in Theorem 2, or
- A sequence of composable control blocks, or
- A control structure whose all subparts are composable control blocks.

1188 **Theorem 4** (Composability of Sequential Control Blocks). *A sequence of composable control*
 1189 *blocks CB_1, CB_2, \dots, CB_n as defined in Definition 9 is composable under the translation function*
 1190 *T .*

1191
 1192 *Proof.* Let CB_1, CB_2, \dots, CB_n be composable control blocks. 1. By Definition 9, each CB_i is
 1193 composable. 2. Translation structure: $T(CB_1 \circ CB_2 \circ \dots \circ CB_n) \equiv T(CB_1) \cdot T(CB_2) \cdot \dots \cdot T(CB_n)$
 1194 where \circ denotes sequential composition in SRC and \cdot denotes concatenation in DST. 3. Control Flow
 1195 Preservation: The sequential order of control blocks is maintained in the translation. 4. Data Flow
 1196 Preservation: The order of operations between control blocks is preserved.

1197 Therefore, the sequence of composable control blocks is itself a composable control block under
 1198 T . \square

1199 **Theorem 5** (Composability of Arbitrary Programs). *Any program P that can be decomposed into*
 1200 *a sequence of control blocks as defined in Definition 4 is composable under the translation function*
 1201 *T if the Iterative Composability Proof algorithm (Figure 2) marks it as composable.*
 1202

1203 *Proof.* The proof follows from the correctness of the Iterative Composability Proof algorithm:

- 1204
 1205 1. The algorithm starts with basic blocks and basic control structures, which are proven com-
 1206 posable by Theorem 1 and Theorem 2.
 1207
 1208 2. It iteratively builds up composability for larger structures:
 1209
 - Sequences of composable blocks are proved composable by Theorem 4.
 - Control structures with all composable subparts are marked composable.
 1210
 1211 3. The process continues until the entire program is marked composable or no further progress
 1212 can be made.
 1213
 1214 4. If the entire program is marked composable, it means that $T(P)$ can be expressed as a
 1215 composition of the translations of its composable parts, preserving both control flow and
 1216 data flow as per Definition 8.
 1217

1218 Therefore, if the algorithm returns that P is composable, then P is indeed composable under the
 1219 translation function T . \square

1220 **Theorem 6** (Composability of Goto Statements). *Goto statements, which are basic statements as*
 1221 *per Definition 2, are composable under the translation function T , but arbitrary goto statements*
 1222 *can break the structured control flow assumed in the main proof.*
 1223

1224 *Proof.* Let l be a label and *goto* l be a goto statement.

1225 Translation structure:

1226
 1227 $T(\text{goto } l) :$
 1228 jmp label_l

1229
 1230 $T(l :) :$
 1231 $\text{label}_l :$
 1232

1233
 1234 The goto statement translates to an unconditional jump, preserving control flow. It doesn't directly
 1235 affect data flow. Composability holds as $T(\text{stmt}_1 \circ \text{goto } l \circ \text{stmt}_2) \equiv T(\text{stmt}_1) \cdot T(\text{goto } l) \cdot$
 1236 $T(\text{stmt}_2)$.

1237 However, goto introduces complications:

- 1238
 1239
 - Non-local control flow can break the nested structure of control blocks.
 - Programs with unrestricted goto usage are difficult to decompose into well-defined control blocks.
 1240
 1241

- It can lead to unstructured code, complicating reasoning about program behavior.

□

While `goto` is provably composable, it's discouraged in modern programming for readability, maintainability, and optimization reasons. Our composability principle is most applicable and valuable in the context of structured programming paradigms.

B.5 SCOPE AND LIMITATIONS OF THE PROOF

The proof of composability presented in this paper is based on a simplified model of C-like languages and unoptimized translation. It's important to note several key points about the scope and limitations of this proof:

1. **Simplification and Correctness:** The simplifications made in our language model and translation process do not compromise the validity of the proof. The core of our argument relies on the decomposition of programs into control blocks and the composability of these blocks. The internal structure of basic blocks, while important for actual compilation, does not affect the composability principle we've established.
2. **Unoptimized Translation:** Our proof assumes a straightforward, unoptimized translation process. This assumption is crucial for maintaining the direct correspondence between source code structures and their translations.
3. **Limitations for Complex Language Features:** The composability principle as proved here can be applied to C-like languages, but may not hold for more complex language features. For example:

- **Exception Handling:** Languages with sophisticated exception handling mechanisms, such as Python, introduce complexities that can break composability. These mechanisms often require:
 - Guarded execution of code blocks.
 - Runtime type information (RTTI) for determining appropriate exception handlers.
 - Non-local control flow that can't be easily decomposed into our model of control blocks.
- **Coroutines and Generators:** Features that allow for suspending and resuming execution mid-function can introduce state that is not easily captured in our model of control flow.
- **Reflection and Metaprogramming:** Languages that allow for runtime modification of program structure or behavior can invalidate static composability assumptions.

Although not applicable to some specific language features, it doesn't mean the composability and its derived LEGO translation method is not applicable to the whole programming language, as long as these features are not used in the code, the composability will still stand and the LEGO translation will still work.

4. **Optimizations Across Basic Blocks:** Our proof assumes that the boundaries of control blocks are respected in the translation process. However, many real-world compiler optimizations operate across these boundaries. Examples include:
 - Loop unrolling
 - Function inlining
 - Global value numbering
 - Code motion optimizations

Such optimizations can reorder, eliminate, or combine operations from different control blocks, potentially breaking the composability property as we've defined it.

5. **Applicability:** Despite these limitations, the composability principle proved here is valuable for:
 - The foundation of LEGO translation method, the proof reveals the composable nature of code in at least control block level, which is a major difference than natural languages.

- The proof process also guided Algorithm 1 in LEGO translation, as proving the composability and making use of the composability share similar algorithms.

In conclusion, while our proof provides a strong foundation for understanding composability in C-like languages with straightforward translation, it's important to recognize its boundaries. More complex language features may require extensions or modifications to this framework to maintain composability guarantees. And optimized code translation usually is not composable.

C DISCUSSIONS

C.1 UNIVERSALITY OF LEGO TRANSLATION

The LEGO translation method, while initially developed for compilation tasks, demonstrates broader applicability based on fundamental properties of programming languages rather than being specific to compilation. The composability that LEGO translation leverages stems from the well-encapsulated control flow and locality principles inherent in modern programming languages (disregarding constructs like `goto` in C, more limitations are clearly described in Appendix E).

These characteristics are intrinsic to programming languages themselves and have guided modern compiler design. They enable the modular partitioning of large-scale programs in modern software development, allowing for incremental and even parallel compilation of code. We harness these properties and apply them to the context of neural compilation using Large Language Models (LLMs).

It's important to note that the applicability of LEGO translation extends beyond compilation. It is suitable for various tasks originating from programming languages, such as code translation between different languages. This method significantly enhances the scalability of machine translation tasks for code, providing a powerful tool for handling large and complex codebases.

C.2 MANAGING HIGHLY COMPLEX EXPRESSIONS

One of the primary challenges in neural compilation arises when dealing with expressions or statements of high complexity. In such cases, LLMs struggle to accurately evaluate these expressions through next token prediction. To address this, we propose two solutions:

- External Tool Integration: We can utilize external parsing tools to generate tree structure information for complex expressions evaluation. This tree structure is then provided to the LLM, offering an explicit traversal order and guiding the evaluation process.
- Expression Decomposition: Without relying on external tools, we can design a new pass where the LLM identifies high-complexity expressions and rewrites them as a combination of lower-complexity expressions. This approach ensures that the entire program consists only of expressions within a proper LLM's evaluation capabilities.

C.3 COMPUTATIONAL COST, EFFECTIVENESS, AND FUTURE PROSPECTS

While our neural compilation method is primarily a proof of concept, it does incur significantly higher computational costs compared to traditional compilation methods - approximately 10^6 to 10^7 times higher. However, this should be weighed against the substantial human resources required for traditional compiler development.

The key advantage of our approach lies in its potential for rapid adaptation to new instruction set extensions or frontend intrinsics. Through techniques like RAG (Retrieval-Augmented Generation) and in-context learning, our method can be extended to support new architectures or language features. This positions neural compilation as a valuable assistant in the compiler development process. A particularly promising application is in generating end-to-end unit tests for compiler adaptation to new instructions. This could significantly streamline the development and testing phases of compiler updates. Recent researches like Munley et al. (2024) have shown the ability to using LLMs to generate unittests during compiler validations.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

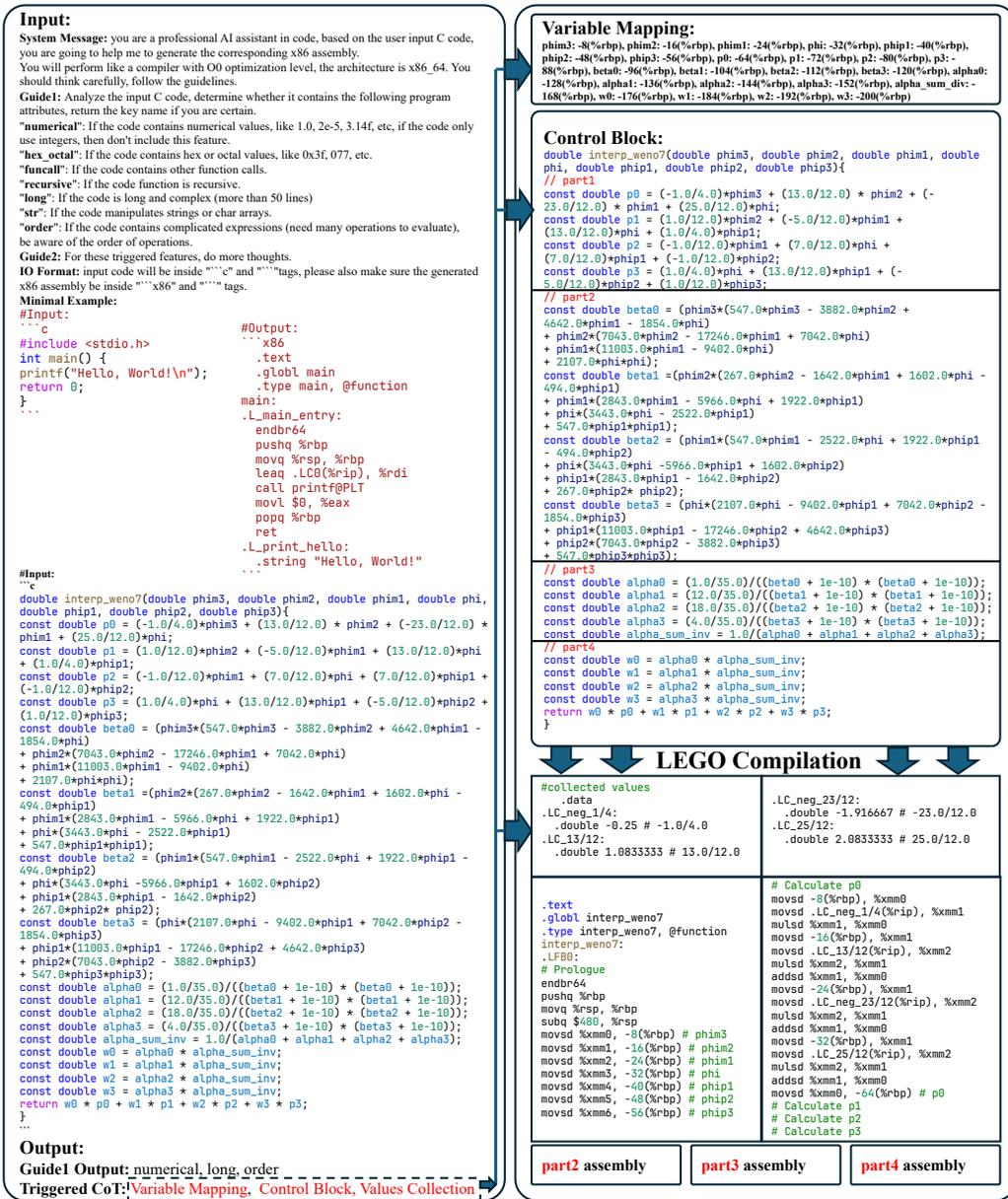


Figure 5: Example workflow for LEGO-Compiler on a full ExeBench example: source code analysis triggers thoughts, including **variable mapping**, splitting **control blocks** and **value collection** illustrated.

D METHOD AND EVALUATION DETAILS

This section provides more details figures, tables and further explanations about **LEGO-Compiler** design and experiment evaluation.

D.1 LEGO-COMPILER: DETAILED DESIGNS

As depicted in Figure 5, LEGO-Compiler is designed to perform a series of thoughts guided by compiler expert knowledge, however, not all CoTs are necessary for each input code, so in our design, we have an **analyze-then-think** approach. First, we will perform an analyzing pass to scan

Table 2: Hardest 10% subset of ExeBench, further breakdown using DeepseekCoder

Ablation1	Baseline	Pass@k	Feedback	CoT	LEGO Translation
DeepseekCoder	63.5%	75.5%	86.0%	96.5%	98.5%
Ablation2	Baseline	CoT	Feedback	LEGO Translation	Pass@k
DeepseekCoder	63.5%	83.5%	90.5%	97.0%	98.5%

the whole program, whose output flags would trigger necessary Chain-of-Thoughts that will be used in the following process. In this example, the code pattern is majorly about double-precision floating point calculations (**numerical**) and complicated expression evaluation (**order**), besides, the code is too long for direct translation method to handle (**long**). Thus, based on the analysis, we applied the following CoTs:

- **Values collection:** A necessary thought, collecting all variables, numericals in a scanning pass, the **numerical** flag will teach the LLM about assembly knowledge to save numerical values.
- **Variable mapping:** Another necessary thought, which will base on the scanned variables and their types, and form a variable mapping table (SymbolTable) for later compilation.
- **Control Block:** the LEGO translation methodology is applied triggered by **long**, where the entire code is considered too long and will be split into control-block level code snippets via Algorithm 1, it’s noteworthy that the **order** flag from analysis will suggest the LLM to split the program into finer-grained blocks so that they can focus more on the order of operations within each block, in Figure 5, there is just one basic block, the flag suggests LLM to split into 4 sequential parts. Then these parts are translated with the aid of SymbolTable individually. Finally, these compiled results are composed together to form a full LEGO compilation.

With different input code, the triggered CoTs will be different, this is helpful because not all thoughts will be useful if no such features appeared in the code, for example, if a code is simple and only has one basic block with a few sequential stmts, then there will be no need to perform LEGO translation related CoTs, because direct translation will be sufficient enough.

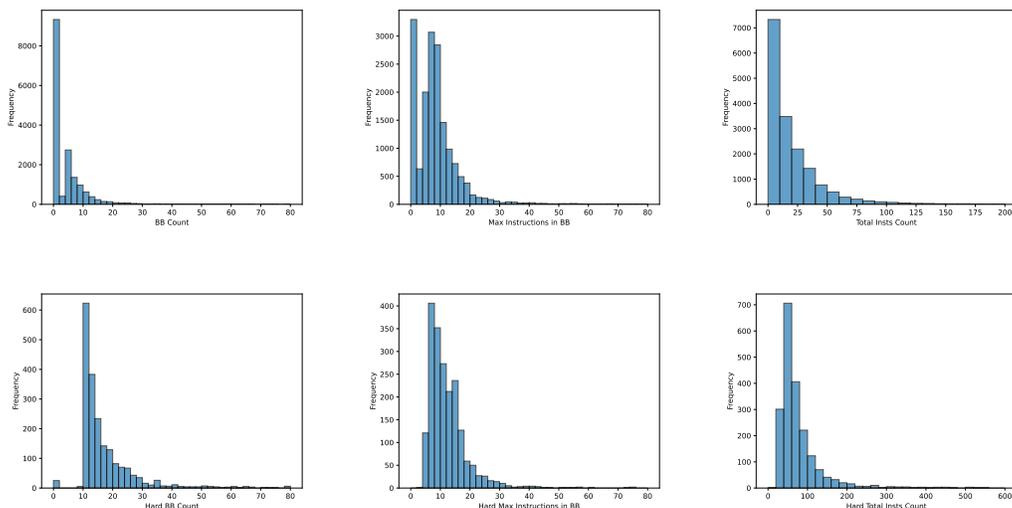
D.2 EXEBENCH BREAKDOWN

Figure 3 shows the complexity ablation on the test set of ExeBench, where we use LLM to categorize all cases into three types of complexity based on certain attributes, like code length, expression complexity, control flow complexity, unusual operations occurrence, etc. The ablation results show LLMs despite of their models’ differences, all get improved on these three categories, where pass@5 and feedback can improve most of the simple cases and some of the medium cases, but can hardly improve on hard cases. While the annotation-based CoT method significantly improve these hard cases, even these mini LLMs can have significant accuracy improvement, except for Codestral model, which fail to follow the CoT correctly, so the result of Codestral for annotation method is a fallback of previous run.

A concern is on whether LLM can categorize code well, so we also perform traditional breakdown, using llvm toolchain (LLVM Project, 2024a) as the frontend analyzer. Based on the analyzed results on basic block count, total instructions and max instructions within a block, we choose the hardest 10% subset of ExeBench for further breakdown. As illustrated in Figure 6, the breakdown characterizes the ExeBench dataset and its hardest 10% subset, which show the subset is significantly harder in total instruction count and basic block count, while the difficulty within each basic block is not significant. After characterization, we use DeepseekCoder as the LLM for evaluation. As depicted in Table 2, although we do find all accuracy degrades due to harder cases, the result further show effectiveness on Feedback, CoT and LEGO translation methods, as the improvement of these methods become more significant.

Furthermore, the hard cases can be majorly categorized into three types:

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476



1477
1478
1479
1480

Figure 6: Complexity breakdown of ExeBench and its hard 10% (roughly) subset, we use llvm as the analysis tool, then filter the subset with the following conditions: number of basic blocks(BB) ≥ 10 or max instructions in BB ≥ 80 or total instructions ≥ 200 . Upper figures characterize the overall of ExeBench and Lower figures characterize the hard 10% subset.

1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502

- The insufficiency on some language-specific features, for example, lacking the knowledge of certain operations, which can be definitely improved with more data in the next model pretrained or by providing external knowledge to aid its generation.
- The unsuccessful reasoning during the annotation-based CoTs. This method require the LLMs to reason arithmetic computation and capture specific code patterns in the code to form intermediate results to aid the generation. If the reasoning process generates incorrectly, the CoTs will fail. However, the reasoning capabilities required for this method is not high, majorly the addition and multiplication of integer values within 1000(typically). As LLMs keep improving their abilities in reasoning and math, this type of failures will reduce significantly.
- Very long code reasoning and follow-up generation, where LLMs fail to generate a very large output at once. The first reason is the limitation of current LLMs themselves, although advanced LLMs have increased their context limits into hundreds of thousands tokens, their single generation capability is still limited, to either 4096, 8192 or 16384 tokens. The second reason is the difficulty to generate a long, error-prone output(like assembly languages) at once, this is an intrinsic drawback of direct generation method itself, and can be solved with the proposed LEGO translation/compilation method. LEGO translation can reduce the complexity to control block level, or at maximum, statement level, however, if the statement itself is very long and complicated to evaluate (which is very rare, but potential in modern programming paradigms), our methods will not help, which is a limitation in our work.

1503
1504

D.3 EVALUATION ON LONGFUNCTION

1505
1506
1507
1508
1509
1510
1511

LongFunction dataset is madeup of 50 C functions in 5 types, where each of them are derived from a certain program pattern like in Figure 7, by alternating the repeated **n**, we could get code size varying from 317 to 238737 tokens, all the token counting is performed by the tiktoken python library, where a gpt-3.5-turbo-0613 vocabulary table is used, although not exactly the token size for each LLM. When evaluating the cases for neural compilation, we compare the neural-compiled results with oracle-compiled results directly since the code is self-contained. As for code translation, we directly test the behavioral output of the translated code and the original code, all compiled by oracle compilers(gcc for C, g++ for C++, CPython runtime for python execution and rustc for

```

1512 #include <stdio.h>
1513 int arr1[10][10];
1514 typedef struct {
1515     float f1;
1516     int i1;
1517 } mystruct1;
1518 typedef struct {
1519     mystruct1 *s1;
1520     int i2;
1521     double d1;
1522     double d2;
1523     mystruct1 *s2;
1524 } mystruct2;
1525
1526 void longfunction1(mystruct2 *res) {
1527     int i, j, k;
1528     // init
1529     for (i = 0; i < 10; i++) {
1530         for (j = 0; j < 10; j++) {
1531             for (k = 0; k < res->i2; k++) {
1532                 arr1[i][j] = 0;
1533             }
1534         }
1535     }
1536
1537     // op1
1538     for (i = 0; i < 10; i++) {
1539         for (j = 0; j < 10; j++) {
1540             for (k = 0; k < res->i2; k++) {
1541                 arr1[i][j] += res->s1->i1;
1542             }
1543         }
1544     }
1545
1546     // op2
1547     for (i = 0; i < 10; i++) {
1548         for (j = 0; j < 10; j++) {
1549             arr1[i][j] += res->s2->i1;
1550         }
1551     }
1552
1553     // op3
1554     for (i = 0; i < 10; i++) {
1555         for (j = 0; j < 10; j++) {
1556             arr1[i][j] += (int)(res->d1 / res->d2);
1557         }
1558     }
1559
1560     // repeat n times(n in 2,4,8,16,32,64,128,256,512)
1561     // will inject a needle in a later repeat pattern
1562     // e.g: += -> -= in op3
1563     return;
1564 }
1565
1566 int main() {
1567     mystruct1 s1 = {1.0, 2};
1568     mystruct1 s2 = {3.0, 4};
1569     mystruct2 result_struct = {&s1, 3, 18.0, 5.0, &s2};
1570     longfunction1(&result_struct);
1571     printf("arr1[0][0] = %d\n", arr1[0][0]);
1572     // other print
1573     return 0;
1574 }

```

Figure 7: LongFunction example code: the code is synthesized by repeating certain patterns with n times, and inject a needle in one of the repeating patterns.

rust). It’s worth noting that the cases in LongFunction are inspired by needle-in-the-haystack experiment (Kuratov et al., 2024), where a needle in the long context must be correctly picked out. In our LongFunction dataset, this is a small, hard to notice modification of the code pattern, for example, replacing a ‘+=’ with ‘-=’. The ability to identify the needle and translate/compile it correctly could significantly support the LLMs with stronger long-context learning ability.

However, if direct translation/compilation is applied, all the models, despite of their long context limits, fail to translate a near 5k token case, and compile a 2.6k token case in LongFunction, and no need to handle all the above. It’s probably LLMs training bias to let it omit similar patterns no matter how we instruct it to step by step thinking and translating.

Our LEGO translation/compilation method, however, can significantly overcome such limitations. Because each time, only a proper sized code snippet is provided to the LLM for further compilation/translation, so theoretically, however long the code is, the LEGO translation method can handle it sufficiently, because small-sized code translation/compilation is assumed to be well-pretrained and proved by results. The splitting and rebuilding processes, although currently not able to be performed all by LLM itself(due to the single output limitation), are simply rule-based and can be well executed by the LEGO-Compiler system, where the splitting process is using the Algorithm 1 algorithm, and the rebuilding process is more simply, concatenating results together.

An easier evaluation can also be performed. By providing an arbitrary code snippet of the long program split following the Algorithm 1, we translate/compile it with the help of globally visible SymbolTable messages and code position markers, if any part of the translation/compilation is semantically correct, then the concatenation of all parts will be correct. This can be easily performed using any LLM api or LLM chat website, and we also provide examples to support this claim.

As a result, all three LLMs (Claude-3.5-sonnet, GPT-4o and DeepseekCoder), successfully translate or compile all the cases in the LongFunction dataset. We also test the capability of newest LLM: o1-preview, although limited to its strict daily usage capacity, it can significantly translate/compile larger sized code snippet, no wonder it can pass all the cases as well.


```

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

```

```

static uint8_t func_1(void)
{
    int64_t l_2[1];
    int32_t l_3 = 0xF37831E4L;
    int32_t l_6[3];
    int i;
    for (i = 0; i < 1; i++)
        l_2[i] = 0xEC2E0CF5720E83C7LL;
    for (i = 0; i < 3; i++)
        l_6[i] = 0xA8CDA2AEL;
    for (l_3 = 0; (l_3 >= 0); l_3 -= 1)
    {
        int16_t l_4 = (-1L);
        int32_t l_5 = (-1L);
        int i;
        l_5 = ((l_2[l_3] != 1UL) <= l_4);
        l_6[0] = l_4;
    }
    l_6[2] = l_3;
    return l_6[0];
}

struct S0 {
    uint8_t f0;
    int32_t f1;
    uint16_t f2;
};

struct S1 {
    struct S0 f0;
    uint32_t f1;
    struct S0 f2;
    uint16_t f3;
};

static struct S1 func_1(void)
{
    uint32_t l_4 = 0xF054A20AL;
    int32_t l_5 = 0x4B03E386L;
    uint8_t l_6[3];
    struct S1 l_11 = {
        {0x8EL, 0x36DC9922L, 0xC436L},
        4294967295UL,
        {1UL, 0xC3FC0233L, 0xD52AL},
        0x2BBDL
    };
    ...
    return l_11;
}

```

Figure 9: Csmith example code, the major body part of the right hand side code is omitted. This example characterizes the necessity of both the Chain-of-Thought reasoning of structs and stack allocation and the LEGO translation method to overcome the complexity of coarse-grained translation.

- The second type of errors is caused possibly by long context forgetting, where the model can not match the current processing assembly with the source code faithfully, LEGO translation method, on the other hand, can handle these cases efficiently as the complexity of each translation is reduced and there are less misleading long contexts to cause these random errors. Besides, finer-grained translation also gives LLMs more attention to faithful translation of operations, the order of operations and implicit conversions.
- The third type is also a limitation our methods can not fully cover: the training bias due to insufficient pretraining in LLMs, which counts for the error in Tripforce’s `generate_password` function, where the translation fails to translate the multiple line strings correctly, which is an insufficient training error in Claude-3.5-Sonnet model itself. Another example is, Claude-3.5-Sonnet model is likely to translate the order of the following expression wrongly: $(x - col * 6)$, when it is a postfix of a lengthy expression, it is likely to generate the subtraction instruction first then the multiplication (causing failures), which is not the case for GPT-4o model and Deepseek model. However, for these models, they have more other training bias that make themselves worse than Claude-3.5-Sonnet model. Using Pass@k and feedback correction can mitigate such failures. Besides, we can be positive about these failures because as LLMs advance, these failures will gradually disappear.

D.5 CSMITH: RANDOMLY GENERATED PROGRAMS EVALUATION

Except for AnsiBench evaluation. We further perform evaluations on randomly generate programs with sufficient complexity. We use Csmith (Yang et al., 2011), a random generator of C programs which is widely used for finding compiler bugs using differential testing as the test oracle. Typically, Csmith examines compilers with random programs with corner case features and numbers, testing the robustness of compilers. Code examples generated from Csmith are illustrated in Figure 9.

As depicted in Figure 10, randomly generated programs by Csmith are very hard for both baseline and CoT-only methods to translate. In a test suite of 25 cases LEGO translation successfully pass, we find baseline translation can only pass 4 cases, with CoT translation, only 9 more cases can be passed. Besides, the complexity of cases only passed by LEGO translation method are signif-

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

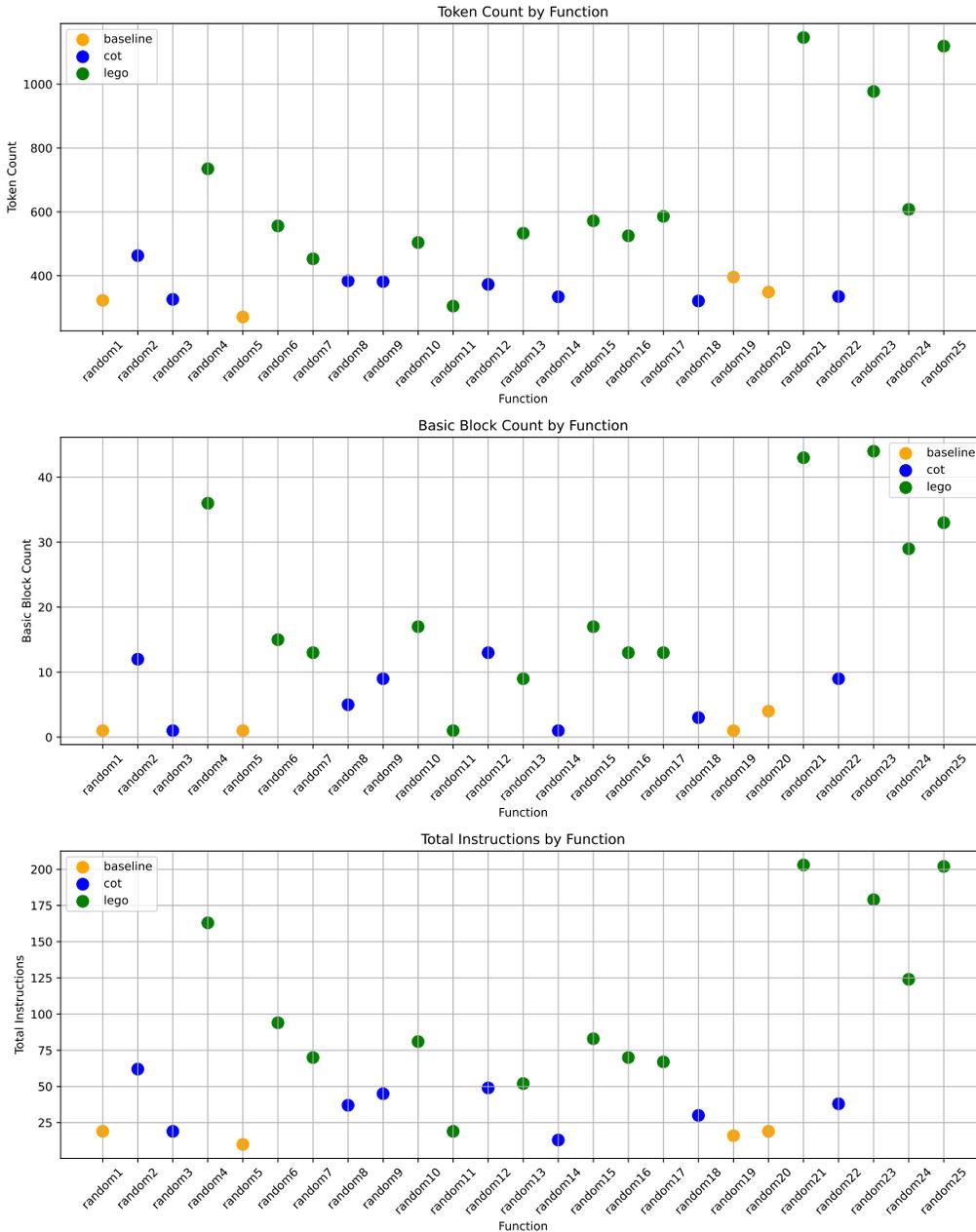


Figure 10: Csmith random generated code statistics, where the practical utility of the LEGO method is shown clearly by passing significantly more complex cases.

Task Description: Analyze every struct in the <code>{src}</code> code, generate the <code>{struct annotation}</code> .	Task Description: Based on <code>{struct annotation}</code> , find all variable definitions in the <code>{src}</code> code, generate the <code>{variable mapping}</code> .	Task Description: Based on <code>{src}</code> code, split the code into middle-sized <code>{control blocks}</code> .	Task Description: Based on <code>{src}</code> code, find all numerical values and literals and save it as <code>{value collection}</code> .
Example	Example	Example	Example
SRC:	SRC:	SRC:	SRC:
<pre>typedef struct { float f1; int i1; } mystruct1;</pre>	<pre>void foo(mystruct2 *res) { int i, j, k; // ... }</pre>	<pre>for (i = 0; i < 10; i++) {...} printf(XXX); // sequential code if(i+j < k){...}</pre>	<pre>const double p0 = (-1.0/4.0)*phi3 +(13.0/12.0)*phi2 +(-23.0/12.0)*phi1 +(25.0/12.0)*phi;</pre>
Struct Annotation:	Variable Mapping:	Control Blocks:	Values Collection:
<pre>typedef struct { float f1;//offset0,size4 int i1;//offset4, size4 } mystruct1;// size 8</pre>	<pre>i:int,[-4,0),-4(%rbp),size 4 j:int,[-8,-4),-8(%rbp),size 4 k:int,[-12,-8),-12(%rbp),size 4 padding [-16,-12) res:ptr,[-24,-16),24(%rbp),size 8</pre>	<pre>// part1 for (i = 0; i < 10; i++) {...} // part2 printf(XXX); // sequential code // part3 if(i+j < k){...}</pre>	<pre>.LC_13: # Literals .double 13.0 # Nums .LC_12: .data .double 12.0 .LC_neg_one: .double -1.0 .LC_neg_23: .double -23.0 .LC_four: .LC_25: .double 4.0 .double 25.0</pre>

Figure 11: Annotation-based Chain-of-Thoughts prompts for neural compilation

Table 3: Ablation study: impact of temperature on Pass@1 and Pass@5 performance

Model	Pass@1					Pass@5				
	0.2	0.4	0.6	0.8	1.0	0.2	0.4	0.6	0.8	1.0
GPT-4o	71%	73%	72%	72%	72%	79%	83%	86%	89%	92%
Claude-3.5-Sonnet	87%	91%	93%	88%	89%	91%	92%	96%	94%	96%
DeepseekCoder	89%	88%	86%	87%	88%	92%	92%	92%	93%	92%
GPT-4o-mini	64%	61%	61%	60%	60%	71%	71%	79%	73%	80%
Claude-3-Haiku	79%	76%	78%	72%	73%	82%	84%	85%	86%	86%
Codestral	73%	66%	41%	-	-	84%	90%	73%	-	-

icantly larger than others, which can be characterized by token count, basic block count and total instructions in the three subfigures respectively.

During Csmith evaluation, we also identify several kinds of errors during LEGO-Compiler translation. For example, overflow value assignment is a kind of error which doesn't usually occur in daily programming but can be found during compiler testing. Taking `int16_t x = 0x56671485;` as an example, it will trigger errors because LLMs directly generate `movw $0x56671485, x'`s address, which fails to check whether the value (overflows the 16 bit word) can be represented through `movw` instruction. Another example is, when handling with implicit type conversions, LLMs may not promote the type correctly, this is critical for floating point computation as operations with wrong precision will cause numerical errors.

D.6 OTHER EVALUATION DETAILS

Table 3 shows the impact of temperature when using LLMs for neural compilation. LLMs have better Pass@1 accuracy when temperature is low, but higher Pass@5 accuracy when temperature is high. This is as expected, since temperature influences the decoding process, with higher temperature, the results are more diverse, allowing LLMs to jump out of pretraining bias, however, this could also cause more errors by choosing sub-optimal decoding tokens that may cause errors.

Figure 11 explains how we prompt LLMs to do the annotation-based Chain-of-Thoughts to aid the neural compilation process.

E LIMITATIONS

Optimization Capabilities: The current focus of LEGO-Compiler is on functional correctness rather than code optimization. Traditional compilers excel at producing highly optimized code, a capability not yet matched by our neural approach. Future work could explore integrating optimization techniques into the neural compilation process.

Performance Overhead: As noted in the discussion, the computational cost of neural compilation is significantly higher than traditional methods. This limitation may restrict its practical application in scenarios where compilation speed is critical.

Complex Expression Handling: The paper acknowledges challenges in managing highly complex expressions, proposing external tool integration or expression decomposition as potential solutions. This indicates a current limitation in LLMs’ ability to handle intricate code structures independently.

Architecture-Specific Knowledge: While the paper demonstrates success with x86, ARM, and RISC-V architectures, expanding to a broader range of architectures, especially more specialized ones, may require significant additional training or fine-tuning of the LLMs, or by providing large RAG database to provide such knowledge in the context.

Security and Reliability: The stochastic nature of LLM outputs raises concerns about the consistency and security of the generated assembly code. Ensuring deterministic outputs and preventing potential vulnerabilities introduced by the neural compilation process remains a challenge.

Handling of Language-Specific Features: The paper primarily focuses on C-like language compilation, and proves the availability of functionality in neural compilation through both theoretical and empirical results. However, extending the approach to other programming languages can result in more tailored problems, for example:

- **RAII idiom:** Languages with class properties, like C++, have an important programming idiom called **Resource Acquisition Is Initialization (RAII)**, which pose significant challenges for LLMs. For instance, constructor and destructor functions in these languages are implicitly called based on scope. This implicit behavior is difficult for LLMs to accurately model and implement in assembly code, but this could be solved using external mangling tools like **c++filt** (Free Software Foundation, 2023).
- **Name Mangling:** Languages like C++ and Rust use name mangling mechanisms for function overloading and template instantiation. This requires special handling of global symbols such as function names during compilation, which may be challenging for LLMs to consistently implement without explicit training on these concepts.
- **Dynamic Language Features:** Some language features violate the composability principle that LEGO translation relies on. For example, Python’s exception handling mechanism, which can cross scope boundaries, would make the LEGO translation method ineffective for such features.

It’s important to note that many of these challenges are not unique to neural compilation. Traditional compilers also struggle with highly dynamic features like exception handling and Run-Time Type Information (RTTI). Languages like Python achieve their flexibility by sacrificing native code generation in favor of interpretation or JIT compilation. Therefore, these limitations are not specific to our work but rather inherent to any approach based on static compilation analysis.

The ability to handle these diverse language features represents an area for future research in neural compilation. It may require developing specialized techniques or combining neural methods with traditional compiler approaches to address these complex language-specific challenges.

Scalability to Very Large Codebases: While the LEGO translation method significantly improves scalability, handling entire large-scale software projects or operating systems may still be beyond the current capabilities of this approach. **However, it is noteworthy that repository complexity is naturally reduced into files or functions, therefore, LLM-based compilers and translators are potential to translate them with more advanced models and more carefully designed methods.**