

# The CodeInverter Suite: Control-Flow and Data-Mapping Augmented Binary Decompilation with LLMs

Peipei Liu\*, Jian Sun\*, Rongkang Sun\*, Li Chen\*, Zhaoteng Yan\*,  
Peizheng Zhang\*, Dapeng Sun\*, Dawei Wang\*, Xiaoling Zhang\*, Dan Li\*†

\*Zhongguancun Laboratory, Beijing, China

{liupp, sunjian, lichen}@zgclab.edu.cn

†Tsinghua University, Beijing, China

tolidan@tsinghua.edu.cn

**Abstract**—Binary decompilation plays a vital role in various cybersecurity and software engineering tasks. Recently, end-to-end decompilation methods powered by large language models (LLMs) have garnered significant attention due to their ability to generate highly readable source code with minimal human intervention. However, existing LLM-based approaches face several critical challenges, including limited capability in reconstructing code structure and logic, low accuracy in data recovery, concerns over data security and privacy, and high computational resource requirements. To address these issues, we develop the CodeInverter Suite, with three main pieces: (1) the CodeInverter Workflow (CIW) is a novel prompt engineering method that incorporates control flow graphs (CFG) and explicit data mappings to improve LLM-based decompilation. (2) Using CIW on well-known source code datasets, we curate the CodeInverter Dataset (CID), a domain-specific dataset containing 8.69 million samples that contains CFGs and data mapping tables. (3) We train the CodeInverter Models (CIMs) on CID, generating two lightweight LLMs (with 1.3B and 6.7B parameters) intended for efficient inference in privacy-sensitive or resource-constrained environments. Extensive experiments on two benchmark datasets demonstrate that CIW significantly enhances the decompilation performance of various LLMs, with average improvements of 9.16% in re-executability and 15.26% in re-compilability. For our proposed decompilation model, CIM-6.7B achieves state-of-the-art performance in terms of re-executability and edit similarity, outperforming existing LLMs—even with over 100× more parameters—by an average of 11.03% and 6.27%, respectively.

## I. INTRODUCTION

Decompilation refers to the process of reversing compiled binary executables back into pcode (pseudo C code) that approximates the original high-level source code [1], [2], [3]. It plays a critical role in various cybersecurity tasks, such as malware analysis[4], vulnerability detection[5], and binary

program auditing[6], [7]. For example, in malware analysis, researchers often have access only to binary executable files due to the unavailability of source code or deliberate obfuscation of developer intent. To understand the logic of the code and uncover potential malicious behavior, it is typically necessary to first decompile the binary into a more readable pcode.

To achieve efficient and accurate decompilation, both the research community and industry have proposed a variety of approaches, which can be broadly categorized into the following four categories: Traditional static analysis tools like Hex-Rays[8] and Ghidra[9], which rely on heuristic rules and schema-matching strategies to recover the basic syntactic and semantic structure of the source code. Machine learning-enhanced methods[10], [11], [3], [12], which integrate learning-based architectures to model the semantics of pcode more effectively. End-to-end decompilation with large language models (LLMs)[13], [14], [15] treat decompilation as a translation task from the assembly code<sup>1</sup> to pcode. Hybrid methods combine LLM with traditional tools[16], [17], [18]. They first use traditional decompilers to generate pcode, and then use LLMs to improve readability and logical completeness.

Among the four, the end-to-end decompilation methods with LLMs are gaining increasing attention in terms of pcode’s readability, syntactic correctness, and reduced manual intervention. However, existing approaches [13], [14], [15] still face three limitations and challenges in practice: (1) They process assembly instructions as plain text sequences for linear modeling, neglecting the inherent structured attributes of the program code. This makes it difficult to capture non-local topological constraints in program execution paths, particularly when handling complex control logic (e.g., nested conditional branches, indirect jumps, and loops), data propagation across basic blocks, and context-dependent relationships, leading to structural fragmentation and logical inconsistencies. (2) They do not consider the data context suggested by the memory

Corresponding author: Li Chen (email: lichen@zgclab.edu.cn).

<sup>1</sup>In binary research works, researchers typically use dedicated tools to first convert binary code into assembly code, which then serves as the foundation for downstream task-specific analyses.

addresses in assembly instruction sequences, hindering the accurate recovery of data objects and their usage relationships during decompilation. For example, in an assembly instruction like `mov eax, [rip+0x3fa29](rip=0x4729)`, the model may treat `[rip+0x3fa29]` merely as a token without considering the actual data information stored at the address `.bss:0x44152`. (3) Most of existing works rely on general-purpose LLMs, including both open-source (e.g., LLaMA[19], StarCoder[20]) and commercial models (e.g., GPT[21], Codex[22]). However, these works face critical challenges in computational resource requirements, data security, and operational costs. For instance, achieving competitive decompilation performance with 70B open-source LLMs requires at least 4x48GB GPUs (with vLLM + FP16) for inference, and utilizing third-party commercial models for decompiling binaries may pose significant data security and privacy risks, especially when the binaries contain proprietary algorithms, business-critical logic, or user data.

To address the above issues, we propose the **CodeInverter Suite**. It consists of three main pieces: the CodeInverter Workflow (CIW), the CodeInverter Dataset (CID), and the CodeInverter Models (CIMs).

The CIW assists LLMs to achieve accurate decompilation. It is a novel prompt engineering method with two core ideas: (1) CIW integrates the semantic modeling strengths of LLMs with the topological insights of control flow graphs (CFGs). By transforming CFGs into sequence-friendly inputs, CIW explicitly incorporates control-flow information into LLM prompts, enabling accurate reconstruction of program structure and control logic. (2) To improve data object recovery, we enhance assembly semantics with an explicit data mapping table extracted from binary data sections such as `.rodata`, `.data`, `.bss`, and `call stack` layouts. The data mapping reflects relationships between raw memory addresses and their corresponding data, and this mechanism makes implicit data dependencies visible to the model and strengthens data object recovery.

Using CIW and public source code datasets, we construct a large, specialized dataset incorporating CFGs and data mapping tables, called CID. Then, to accommodate privacy-sensitive and resource-constrained deployment environments, we use CID to train two decompilation LLMs with 1.3B and 6.7B parameters (CIM-1.3B/6.7B). Compared to prior work, our compact models not only offer several distinct advantages such as deployability, lower computational resource requirements, and complete security control but also significantly improve performance. Experimental results show our models achieve superior performance across major decompilation benchmarks (re-executability, Pass@k, and readability).

In summary, our contributions in this paper can be summarized as follows:

- We design the CIW, which enables the integration of control flow graphs and data mapping tables into LLM inputs to enable end-to-end decompilation. This approach leverages control flow information to reconstruct program structure and uses data context to recover data objects. To

the best of our knowledge, this is the first work to incorporate both information into LLMs for decompilation.

- We construct the CID and have trained two models (CIM-1.3B/6.7B) based on this dataset. The dataset includes 8.69 million assembly-source code pairs for decompilation, along with semantic labels aligned with control flow graphs and data mapping tables for memory addresses. Our models provide a practical solution for end-to-end decompilation in resource-constrained environments.
- We demonstrate the effectiveness of our proposed method and the superiority of our decompilation LLMs through extensive experiments and analysis across various baseline LLMs.
- We open-source our CID and CIMs to facilitate community research and applications at <https://github.com/LiuPeiP-CS/CodeInverter>.

## II. BACKGROUND & MOTIVATION

### A. Control Flow Graph

CFG is a widely adopted intermediate representation in binary program analysis. It decomposes programs into basic blocks and explicitly models transfer relationships between them (e.g., conditional jumps, loops, and function calls), thereby characterizing all potential execution paths in a program. Compared to linear representations, CFGs provide more accurate topological structures and control dependency information from a static perspective, enabling models to better understand control constraints and dependency relationships between different code blocks.

The application of CFGs has been widely explored in binary research works. CodeCMR [23] performs binary-source code matching by learning structural program features over CFGs using graph neural networks (GNNs). Order Matters [24] employs convolutional neural networks on adjacency matrices of CFGs to extract topological ordering features for binary code similarity analysis. NeurDP [10] employs GNNs over CFGs to address challenges posed by compiler-optimized binaries, introducing intermediate representations and Optimized Translation Units to improve the decompilation of optimized binary code. In the type inference task of TYGR [11], the CFG provides the structural foundation for exploring functions, guiding the traversal of program paths and enabling precise data-flow analysis. In the work of [3], the CFG is leveraged to encode the structural context of call sites within procedures, enhancing semantic understanding and guiding neural models in procedure name prediction. NFRE [25] introduces a structure-sensitive instruction embedding method based on CFGs to support function name prediction.

### B. Decompilation with LLM

Recent advancements in decompilation have been motivated by the success of LLMs. DecGPT [16], DeGPT [17] and WaDec[18] employ LLMs to refine and enhance the output of traditional decompilation tools.

```

1 static void setup_nhm32(void)
2 {
3     static float possible_nhm_bus[] = {0xFF, 0x7F, 0x3F};
4     unsigned long did, vid, mc_control, mc_ssrcontrol;
5     int i;
6     ctrl.cap = (2 | 4 | 8);
7     ctrl.mode = 0;
8     for (i = 0; i < sizeof(possible_nhm_bus) / sizeof(
9         possible_nhm_bus[0]); i++) {
10         pci_conf_read(possible_nhm_bus[i], 3, 4, 0x00, 2, &vid);
11         pci_conf_read(possible_nhm_bus[i], 3, 4, 0x02, 2, &did);
12         vid ^= 0xFFFF;
13         did ^= 0xFF00;
14         if (vid == 0x8086 && did >= 0x2C00) {
15             nhm_bus = possible_nhm_bus[i];
16         }
17     }
18     pci_conf_read(nhm_bus, 3, 0, 0x48, 2, &mc_control);
19     if ((mc_control >> 1) & 1) {
20         ctrl.mode = (2 | 4);
21         pci_conf_read(nhm_bus, 3, 2, 0x48, 2, &mc_ssrcontrol);
22         if (mc_ssrcontrol & 1) {
23             ctrl.mode = (2 | 4 | 8);
24         }
25 }

```

(a) Source code

```

1 void setup_nhm32()
2 {
3     unsigned int did;
4     unsigned int vid;
5     unsigned int mc_control;
6     unsigned int mc_ssrcontrol;
7     int i;
8     unsigned int v5;
9     v5 = __readgsdword(0x14u);
10    ctrl.cap = 14;
11    ctrl.mode = 0;
12    for (i = 0; (unsigned int)i <= 2; ++i) {
13        pci_conf_read(*(float *) (char *) setup_nhm32::possible_nhm_bus + (_DWORD)&
14            dword_0[i]), 3, 4, 0, 2, &vid);
15        pci_conf_read(*(float *) (char *) setup_nhm32::possible_nhm_bus + (_DWORD)&
16            dword_0[i]), 3, 4, 2, 2, &did);
17        vid = (unsigned __int16)vid;
18        did ^= 0xFF00u;
19        if ( (unsigned __int16)vid == 32902 && did > 0x2BFF )
20            nhm_bus = *(float *) (char *) setup_nhm32::possible_nhm_bus + (_DWORD)&
21                dword_0[i]);
22    }
23    pci_conf_read(nhm_bus, 3, 0, 72, 2, &mc_control);
24    if ( (mc_control & 2) != 0 ) {
25        ctrl.mode = 6;
26        pci_conf_read(nhm_bus, 3, 2, 72, 2, &mc_ssrcontrol);
27        if ( (mc_ssrcontrol & 1) != 0 )
28            ctrl.mode = 14;
29    }
30 }

```

(b) Hex-Rays

```

1 void setup_nhm32(void)
2 {
3     static float possible_nhm_bus[] = { 255.0, 127.0, 63.0 };
4     unsigned int i, dev, ven, mc_control, mc_ssrcontrol;
5     ctrl.cap = 0xe;
6     ctrl.mode = 0;
7     for (i = 0; i < sizeof(possible_nhm_bus) / sizeof(
8         possible_nhm_bus[0]); i++) {
9         pci_conf_read(possible_nhm_bus[i], 3, 4, 0, 2, &dev);
10        pci_conf_read(possible_nhm_bus[i], 3, 4, 2, 2, &ven);
11        dev ^= 0xffff;
12        ven ^= 0xffff;
13        if (dev == 0x8086 && ven >= 0x2c00) {
14            nhm_bus = possible_nhm_bus[i];
15            i++;
16            break;
17        }
18    }
19    pci_conf_read(nhm_bus, 3, 0, 72, 2, &mc_control);
20    if (mc_control & 0x2) {
21        ctrl.mode = 0x6;
22        pci_conf_read(nhm_bus, 3, 2, 72, 2, &mc_ssrcontrol);
23        if (mc_ssrcontrol & 0x1) {
24            ctrl.mode = 0xe;
25        }
26    }
27 }

```

(c) CIM w/o CFG

```

1 void setup_nhm32(void)
2 {
3     static float possible_nhm_bus[] = {
4         0x00000000, 0x00000001, 0x00000002, ..., 0x000000fe, 0x000000ff,
5     }; //For display convenience, some values are omitted.
6     static float nhm_bus;
7     unsigned int dev_id, ven_id;
8     unsigned int mc_control, mc_ssrcontrol;
9     ctrl.cap = 0x0000000e;
10    ctrl.mode = 0x0000000e;
11    for (i = 0; i < sizeof(possible_nhm_bus) / sizeof(float); i++) {
12        if (i < 3) {
13            pci_conf_read(possible_nhm_bus[i], 0x03, 0x04, 2, &dev_id);
14            pci_conf_read(possible_nhm_bus[i], 0x03, 0x02, 2, &ven_id);
15            dev_id ^= 0xffff;
16            ven_id ^= 0xffff;
17            if ((ven_id == 0x8086) && (dev_id >= 0x2c00))
18                nhm_bus = possible_nhm_bus[i];
19        }
20    }
21    pci_conf_read(nhm_bus, 0x00, 0x48, 2, &mc_control);
22    if (mc_control & 0x0002) {
23        ctrl.mode = 0x00000006;
24        pci_conf_read(nhm_bus, 0x03, 0x48, 2, &mc_ssrcontrol);
25        if (mc_ssrcontrol & 0x0001)
26            ctrl.mode = 0x0000000e;
27    }
28 }

```

(d) CIM w/o data mapping

```

1 void setup_nhm32(void)
2 {
3     static float possible_nhm_bus[] = { 255, 127, 63 };
4     unsigned int i, dev, ven, mc_control, mc_ssrcontrol;
5     ctrl.cap = 0xe;
6     ctrl.mode = 0;
7     for (i = 0; i < sizeof(possible_nhm_bus)/sizeof(
8         possible_nhm_bus[0]); i++) {
9         pci_conf_read(possible_nhm_bus[i], 3, 4, 0, 2, &dev);
10        pci_conf_read(possible_nhm_bus[i], 3, 4, 2, 2, &ven);
11        dev ^= 0xffff;
12        ven ^= 0xffff;
13        if (dev == 0x8086) && (ven >= 0x2c00) {
14            nhm_bus = possible_nhm_bus[i];
15        }
16    }
17    pci_conf_read(nhm_bus, 3, 0, 0x48, 2, &mc_control);
18    if (mc_control & 0x2) {
19        ctrl.mode = 0x6;
20        pci_conf_read(nhm_bus, 3, 2, 0x48, 2, &mc_ssrcontrol);
21        if (mc_ssrcontrol & 0x1) {
22            ctrl.mode = 0xe;
23        }
24    }
25 }

```

(e) CIM

Fig. 1: Motivating example. Presented are the source code (a) of the case sample alongside the decompilations of Hex-Rays (b), CIM without CFG (c), CIM without data mapping (d) and CIM (e).

Unlike the aforementioned approaches, other research efforts have constructed LLMs to perform direct assembly-to-source decompilation. Nova [14] uses hierarchical attention and contrastive learning to enhance accuracy. Feng et al. [15] proposes FAE, which recompiles decompiled code for in-context learning and better alignment with source code. LLM4Decompile [13] supports both direct binary decompilation and Ghidra output refinement, achieving strong results in readability and executability.

### C. Motivations

Given the proven utility of structural features derived from CFGs in binary research works, and the increasing adoption of LLMs in decompilation tasks, we aim to explore the integration of CFGs with LLMs to improve the intelligence and effectiveness of the decompilation process. Beyond control flow, we further observe that explicitly mapping memory addresses to their corresponding data contexts can significantly enhance the recovery of data objects, especially in complex binary programs. Such mappings help reveal the relationships between variables, their types, and usage patterns—information that is often obscured or lost during compilation.

To clearly illustrate the motivation behind our approach, we present a concrete example that demonstrates the benefits

of incorporating both control-flow structures and data-related contextual information into an LLM-based framework.

The motivating example is the function `setup_nhm32` from the ExeBench [26] dataset, whose source code is shown in Fig.1(a). The function’s behavior relies on two critical elements: (1) the initialization of the variable `possible_nhm_bus` (line 3), which is propagated to its caller functions, and (2) the `if` branch at line 13, where computations influence subsequent control flow. When compiled with `-O0` for the x86-32 architecture and decompiled using the Hex-Rays[8] decompiler (Fig. 1(b)), we observe compound failures: the variable `possible_nhm_bus` remains uninitialized, and `vid` is incorrectly assigned. These issues jointly corrupt the subsequent branch execution, illustrating the interdependence between value recovery and control-flow reconstruction—errors in one domain can propagate and compromise the other.

These limitations motivate the design of CIMs. The data-mapping-only variant accurately recovers constants (Fig. 1(c)) but produces incoherent control structures (e.g., lines 14–15). In contrast, the CFG-only version (Fig. 1(d)) preserves the logical control flow but severely misrecovers values—for instance, expanding the variable `possible_nhm_bus` to 256 data entries, whereas the source defines only 3. These

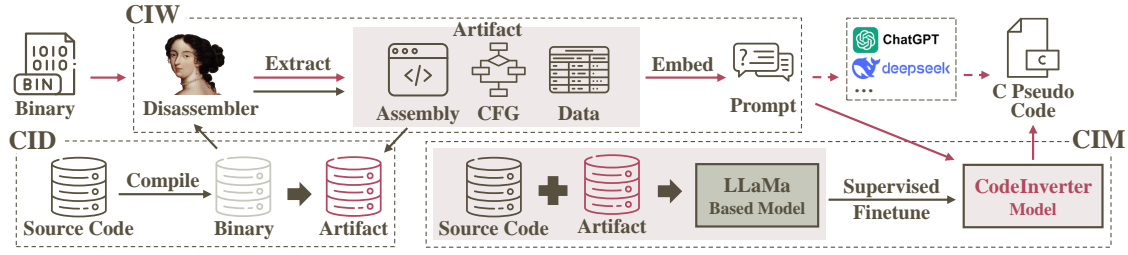


Fig. 2: The decompilation workflow with LLMs in this paper (top: inference, bottom: training).

errors cascade: incorrect values distort the behavior of caller functions, and flawed control flow disrupts value propagation paths. Only the full CIMs (Fig. 1(e)), which integrates both components, achieves functionally correct decompilation. While minor representational differences remain, the synergy of CFG-guided structural accuracy and precise value recovery resolves the interdependence problem observed in Hex-Rays. This demonstrates that reliable decompilation requires the joint validation of macro-level control structures and micro-level data values—a capability uniquely enabled by our work.

### III. DESIGN OF CODEINVERTER SUITE

#### A. CFG-augmented Code Structure Reconstruction

In decompilation tasks, achieving code structure and control logic consistent with the original source code requires LLMs to effectively analyze and leverage features including branch logic, loop boundaries, and functional hierarchies in assembly instructions. However, raw assembly instructions are presented as linear sequences and lack explicit representations of these features, making it difficult for LLMs to accurately recover the source code. Considering that the CFG can explicitly model a program’s execution logic and provide LLMs with structured topological cues, we propose to integrate CFGs with LLMs to facilitate the code structure reconstruction.

**CFG Extraction.** Based on the assembly code obtained from the input binary, we extract the CFG for target functions using the IDA disassembler [27], with assistance from objdump [28] to improve function boundary identification.

The process begins with function identification: We first use IDA to determine function boundaries. When symbol tables are available, we identify function entry points by matching function symbols with their corresponding addresses in the binary. In cases where IDA fails to identify a target function, we use the starting address derived from objdump to specify the entry point for subsequent control flow analysis.

After that, basic block segmentation is performed within each function recursively, where the instruction sequence is split at each control transfer instruction (e.g., jumps, calls) and their immediate targets. Then, edges between basic blocks are established through direct branch resolution for unconditional and conditional jumps, and through fall-through analysis for sequential execution paths.

**CFG Representation.** The extracted CFG is first represented as a structured dictionary incorporating well-labeled assembly instruction blocks (as shown in Figure 3), and then

the dictionary is serialized into a JSON string to meet the input requirements of LLMs.

To improve both machine understanding and processing efficiency, we design a semantically enriched symbolic labeling scheme. Each basic block is assigned a unique symbolic label, and the target transfer addresses in control flow instructions are replaced with the corresponding block labels. For instance, the instruction `jmp short 0x473a` is rewritten as `jmp short loc_473A`, where `loc_473A` is the label assigned to the target block. Similarly, function call operands are replaced with function identifiers instead of numeric addresses, such as `puts_bench` in `loc_473A` block.

Subsequently, we represent the CFG in a dictionary (“Control Flow Graph” in Figure 3) consisting of three key fields: “nodenum” (the total number of basic blocks), “nodes” (the instruction content of each block), and “edges” (a list of directed edges between blocks). Each edge is represented as a pair of source and target block labels—for example, [“start”, “loc\_4729”] indicates a control transfer from the `start` block to the `loc_4729` block.

#### B. Data Mapping-augmented Data Object Recovery

During the program compilation, most data objects—including variables and constants—are stored in designated non-code memory regions, with the exception of a small subset of constants that are embedded directly as immediate values within the code region `.text`. These memory regions are organized as follows: the `.rodata` section stores immutable data known at compile time, including string literals, constant arrays, and other const-qualified variables; the `.data` section contains global and static variables that are explicitly initialized with non-zero values; the `.bss` section holds global and static variables that are uninitialized in the source code, which are zero-initialized at runtime; the `call stack` holds temporary variables (e.g., local variables, function arguments) and return addresses, facilitating runtime context preservation and function invocation.

In binary decompilation, distinguishing between the `.data` and `.bss` sections aids in inferring variable types and recovering global/static data objects. The `.rodata` section contains both semantically rich data—such as string literals—and structural elements like jump tables, which are critical for reconstructing constant references and control-related data objects, respectively.

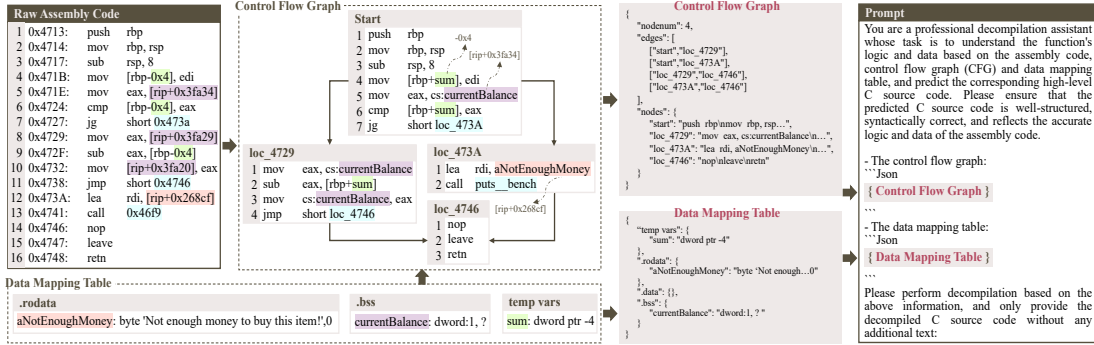


Fig. 3: Prompt engineering details of our proposed CIW (the resulting prompt for all LLMs)

Therefore, leveraging these data sections can improve the completeness and accuracy of source-level data object recovery compared to prior `.text`-only methods. Although ReF Decompile [29] considers such sections, it relies on metadata extracted from the source code, making it impractical and unrealistic in real-world scenarios where the source code is unavailable. In contrast, we provide a more practical and principled solution by exploiting cues from the binary-level data sections.

**Data Retrieval.** We build an IDA plugin to retrieve the data associated with each target function. The process begins by identifying relevant data sections (e.g., `.rodata`, `.bss`) through parsing the section headers. These sections define the scope of static memory regions that may contain data used by the program.

After identifying these sections, we extract their full contents to form a pool of candidate data items for further analysis. We then perform cross-reference analysis to determine which data items are validly referenced in the assembly instructions. Data items that are explicitly referenced within the address range of the target function are retained as directly associated. For data items without explicit references, we additionally include those that are in close proximity to already retained items, under the assumption that such proximity may indicate memory adjacency (e.g., array elements). This heuristic helps reduce the risk of omitting relevant data that may be implicitly referenced.

Additionally, we identify temporary stack variables by analyzing stack register addressing and offset values, as they are crucial for recovering the types and semantics of underlying data objects.

**Data Mapping Table.** With the retrieved data, we subsequently replace the raw address identifiers with semantically rich labels, such as `[rip+0x3fa34](rip=0x471E)→.bss:0x44152→currentBalance` and `-0x4→sum` (Figure 3). As a result, we obtain a data mapping table corresponds to the target function, assigning meaningful data items to the memory addresses in the assembly instructions.

As shown in Figure 3, each data item includes its size and corresponding value. Specifically, data sizes are denoted using standard types: `byte`, `word`, `dword` (double word), and `qword` (quad word); values are primarily represented in

hexadecimal format, except when strings are detected—such as “Not enough money to buy this item!” Strings are identified by scanning for sequences of printable ASCII or Unicode characters terminated by null bytes. Based on contextual and section attributes, appropriate string types (e.g., C-style or Pascal-style) are then assigned. Uninitialized items in the `.bss` section are denoted with a “?”. To concisely represent consecutive uninitialized items of the same size in the `.bss` section—e.g., array-like structures—a count is appended to the data size. For example, `dword:1, ?` indicates a single uninitialized `dword`-sized value. For stack variables, only the offset and data size are recorded, such as `temp vars` in Figure 3.

Similar to the representation of the CFG, the data mapping table is structured as a dictionary and then serialized into a JSON string, allowing LLMs to effectively interpret and utilize the information.

### C. CID Dataset Construction

We construct our dataset based on the existing source code dataset ExeBench [26], a function-level dataset that pairs real-world C functions from GitHub with input-output (IO) examples to enable executable supervision. Starting with 1.2 million C functions from ExeBench, we generate our dataset with 32-bit and 64-bit through the following steps (as illustrated in the bottom left of Figure 2):

**Step1: Resolve Data Type Conflicts.** To generate 32-bit executables of C functions, we address data type conflicts between 32-bit standard libraries and ExeBench-generated `.c` files. For instance, we replace the 64-bit-specific definition `typedef unsigned long size_t;` in ExeBench with `typedef unsigned int size_t;` to comply with the 32-bit program standards.

**Step2: Compile.** The `.c` files containing the source code of the target functions are compiled using GCC into executables with four different optimization levels (O0, O1, O2, and O3). After compilation, the executables are stripped to remove debugging information for further processing.

**Step3: Generate Disassembly Artifact.** We disassemble the binary executables using IDA, and extract the corresponding CFGs and data mapping tables based on the processing techniques of CIW (i.e., Section III-A and Section III-B).

**Step4: Data Representation.** We design a decompilation-oriented LLM prompt and integrate the CFG and data mapping table into it (as shown at the right of Figure 3).

#### D. CIM Training and Inference

We adopt a sequence-to-sequence (seq2seq) framework to train our decompilation LLM, which maps assembly instruction sequences to corresponding source code sequences. We initialize the model with checkpoints from LLM4Decompile [13], which is based on the LLaMA architecture [19]. The training objective follows the standard approach used in neural machine translation, which is well-suited for decompilation tasks due to their sequential generation nature: given an input sequence, the model learns to predict the corresponding output sequence token by token.

**Output Prediction Probability.** Given an input assembly sequence  $\mathbf{x}_{\text{in}} = [x_1, \dots, x_m]$  and a target sequence of source code  $\mathbf{x}_{\text{out}} = [x_{m+1}, \dots, x_n]$ , the model predicts the probability of  $i$ -th token autoregressively using a softmax distribution over the decoder’s output logits:

$$P(x_i | x_{1:i-1}; \theta) = \text{softmax}(\mathbf{W}_o \mathbf{h}_i^{\text{dec}})_{x_i}, \quad (1)$$

where  $\mathbf{h}_i^{\text{dec}}$  is the decoder’s hidden state at position  $i$ ,  $\mathbf{W}_o$  is the output projection matrix, and  $\theta$  denotes all trainable parameters. During training, we apply teacher forcing by conditioning the prediction of  $x_i$  on the ground-truth prefix  $x_{1:i-1}$ .

**Loss Function.** The model is trained by minimizing the cross-entropy loss between the predicted token distributions and the ground-truth source code tokens. For a sequence of total length  $n$ , the loss is defined as:

$$\mathcal{L}(\theta) = - \sum_{i=m+1}^n \log P(x_i | x_{1:i-1}; \theta), \quad (2)$$

where the summation starts at  $i = m + 1$  to ensure that only the target source code tokens contribute to the loss. This is equivalent to maximizing the conditional likelihood of the target sequence  $\mathbf{x}_{\text{out}}$  given the input sequence  $\mathbf{x}_{\text{in}}$  under the model parameters  $\theta$ .

**Inference Process.** During inference, given a disassembled input sequence  $\mathbf{x}_{\text{in}} = [x_1, \dots, x_m]$ , the model generates the target source code tokens autoregressively. At each decoding step  $i > m$ , the model predicts the next token based on the entire input sequence and previously generated tokens:

$$\hat{x}_i = \arg \max_x P(x | x_1, \dots, x_m, \hat{x}_{m+1}, \dots, \hat{x}_{i-1}; \theta), \quad (3)$$

where  $\hat{x}_{m+1}, \dots, \hat{x}_{i-1}$  denote the tokens generated so far.

Generation proceeds iteratively until an end-of-sequence token is produced or a maximum length constraint is reached. Beam search or sampling-based decoding strategies (e.g., nucleus sampling, top- $k$  sampling) can be employed to enhance generation quality.

## IV. EXPERIMENTS

### A. Experimental Setups

**Evaluation Datasets.** Following prior works [13], [15], [14], we use HumanEval-Decompile [13] and ExeBench test set [26] as evaluation datasets.

*HumanEval-Decompile* consists of 164 C functions derived from Python solutions and assertions in HumanEval [30]. Each function is compiled using `gcc` with standard C libraries and passes all assertions. All functions are compiled with optimization levels ranging from `-O0` to `-O3`, then disassembled into `x86_64` assembly before being used for decompilation. HumanEval-Decompile is a widely used benchmark for evaluating code generation. Since the original HumanEval-Decompile only supports `x86_64` compilation, we extend the dataset to support 32-bit compilation for generalization testing. For the sake of experimental design convenience, we abbreviate HumanEval-Decompile as HumanEval in this paper.

In contrast, the original test set of *ExeBench* contains 5,000 real-world C functions extracted from GitHub repositories. These programs include not only complete function definitions but also input-output (I/O) examples and corresponding external functions or header files, ensuring that each function is executable. Compared to HumanEval-Decompile, many functions in ExeBench involve user-defined data structures and more complex dependencies. We also extend the dataset to support 32-bit compilation.

**Evaluation Metrics.** Following prior work, we adopt Re-compilation Rate (Re-com), Re-executability Rate (Re-exe), Edit Similarity (ES), and Pass@ $k$  as evaluation metrics to assess decompilation performance [13], [31], [14], [32].

*Re-compilation Rate:* This metric evaluates whether the decompiled code generated by our model can be successfully recompiled into an executable binary without errors. A high Recompilation Rate indicates that the decompiled code is syntactically correct and adheres to the constraints of the target programming language (e.g., C).

*Re-executability Rate:* This metric evaluates the functional correctness of the decompiled code. It measures whether the recompiled binary can execute successfully and produce outputs that accurately match the expected results. A high Re-executability indicates that the decompiled code is semantically and functionally correct, preserving the original program’s logical behavior.

The *Edit Similarity* is a key metric for evaluating the readability of decompiled code, and it quantifies the similarity between the decompiled code and original code. Following the work of [31], we use edit distance, a standard metric employed in other neural approaches [33], [34], to define edit similarity. The result,  $ES(A, B)$ , is normalized by the length of the ground truth sequence, where a higher edit similarity indicates better readability of the decompiled code. Specifically, it is calculated as follows:

$$ES(A, B) = 1 - \frac{ED(A, B)}{L_B} \quad (4)$$

where  $A$  represents the prediction sequence,  $B$  represents the true sequence,  $L$  represents the length of the sequence, and  $ED(A, B)$  represents the editing distance between  $A$  and  $B$ . Editing distance is also called Levenshtein distance, which means the minimum number of operations required to convert a sequence  $A$  to  $B$ . It is usually computed by dynamic programming, and its state transition equation is as follows:

$$ED(A_i, B_j) = \begin{cases} \max(L_A, L_B) & \min(L_A, L_B) = 0 \\ \min \begin{cases} ED(A_{i-1}, B_j) + 1 \\ ED(A_i, B_{j-1}) + 1 \\ ED(A_{i-1}, B_{j-1}) + 1 \end{cases} & \text{otherwise} \end{cases} \quad (5)$$

**Pass@ $k$  Metric:** In the decompilation task, Pass@ $k$  measures the probability that, given  $n$  decompiled outputs generated by the model for a single binary function ( $n \geq k$ ), at least one of the  $k$  randomly selected outputs passes a functional correctness check (i.e., correct execution against test cases). This metric is widely used to evaluate the success rate of models in scenarios where multiple candidate outputs are allowed.

To estimate Pass@ $k$ , we adopt the unbiased estimator proposed in [32], defined as:

$$\text{pass@}k := \mathbb{E}_{\text{Problems}} \left[ 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right], \quad (6)$$

where  $n$  is the total number of generated outputs,  $c$  is the number of correct outputs (i.e., those that pass execution or functional validation), and  $k \leq n$  is the number of samples randomly selected.

This estimator computes the expected probability that at least one correct sample is included among the  $k$  draws. It provides a balanced view of the model’s ability to produce diverse yet accurate outputs, which is crucial for decompilation tasks where a single deterministic output may not capture all possible valid reconstructions.

For instance, suppose the model generates  $n = 20$  candidate decompilation outputs for a given binary function, and only one of them passes the correctness check. Then, the empirical estimate of pass@1 is  $1/20 = 0.05$ .

Compared to the naive estimator  $1 - (1 - \hat{p})^k$  (where  $\hat{p}$  is the empirical pass@1), Equation 6 provides an unbiased estimate and avoids overestimating performance when  $c$  is small. In our implementation, we follow a numerically stable computation strategy as suggested in [32] to handle large combinatorial terms efficiently.

**Baseline LLMs.** To demonstrate the effectiveness of our proposed **CIW** and the superiority of **CIM**, we conduct both vertical and horizontal evaluations. The vertical evaluation is designed to highlight the importance of **CIW** in the decompilation process, while the horizontal evaluation compares the decompilation performance of our LLMs with other existing LLMs.

Based on this setting, we select general-purpose LLMs with strong generalization abilities (e.g., GPT-4o [35], DeepSeek-V3 [36], and Qwen-plus [37]), along with decompilation-

specific LLMs (e.g., LLM4Decompile [13], FAE [15], and Nova [14]) as baselines.

- LLM4Decompile [13], the latest and most advanced LLM for end-to-end decompilation, built on DeepSeek-Coder. It supports both direct binary decompilation and the refinement of Ghidra’s outputs, delivering significant improvements in readability and executability. In this paper, we compare with its direct binary decompilation version.
- GPT-4o [35], a top-performing general-purpose LLM, renowned for its outstanding capabilities in understanding and generating high-level programming languages.
- DeepSeek-V3 [36], an advanced open-source language model designed for high-performance natural language processing tasks, optimized for multilingual understanding and generation. It emphasizes efficiency and scalability, supporting diverse applications like code generation, mathematical reasoning, and context-aware dialogue. The version we use is DeepSeek-V3 671B.
- FAE [15] utilizes debugging information to accurately align assembly code with source code at the statement level. And it is a 6.7B model in the work of [15].
- Qwen-plus [37], a high-performance proprietary LLM optimized for complex language understanding and generation tasks. Building upon the foundation of its predecessor, Qwen, it delivers enhanced scale, capabilities, and efficiency. Its inference performance, cost, and speed are positioned between Qwen-Max and Qwen-Turbo, making it ideal for moderately complex tasks. The version currently in use is Qwen-plus-2025-01-25.
- Nova [14] also introduces two LLMs (1.3B and 6.7B) with hierarchical attention and contrastive learning objectives to improve the accuracy of decompilation. Nova achieves SOTA decompilation performance on the Pass@ $k$  metric.

**Training and Inference.** We build upon LLM4Decompile [13] as the base model and fine-tune both 1.3B and 6.7B parameter variants for one epoch (15 steps) using the LoRA technique [38], with a rank of 32 and an alpha of 64. The fine-tuned modules include the embedding layer, the language model head, and all projection layers. The batch sizes are set to 16 and 12 for the 1.3B and 6.7B models, respectively. We employ the AdamW optimizer with an initial learning rate of  $2e-5$ , and set the maximum sequence length to 4096. Training is conducted on an NVIDIA GPU cluster with 8×H100-80GB, and we utilize ColossalAI for efficient distributed training.

To enhance generation performance, we adopt a hierarchical training strategy, gradually introducing samples from simple to complex based on the number of CFG blocks. During inference and evaluation, we consistently apply greedy decoding. The decompilation process is accelerated using vLLM [39], which also ensures deterministic output.

For the Pass@ $k$  metric, we follow the same inference configuration (temperature = 0.2, top\_p = 0.95) as in Nova [14], and report Pass@1 and Pass@10 on the HumanEval 64-

TABLE I: Decompile results on HumanEval with only assembly instruction (without **CIW**). (%)

Metric	Model	HumanEval 32-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o	85.98	82.32	83.54	78.05	82.47	89.02	77.44	85.98	79.27	82.93
	Deepseek-V3	<b>94.51</b>	<b>87.80</b>	<b>85.98</b>	<b>89.63</b>	<b>89.48</b>	<b>95.73</b>	<b>88.41</b>	<b>89.02</b>	<b>84.15</b>	<b>89.33</b>
	Qwen-plus	25.61	37.80	41.46	46.34	37.80	73.17	58.54	61.69	62.20	63.90
Re-exe	GPT-4o	22.56	11.59	13.41	9.51	14.27	33.54	10.98	14.02	9.76	17.08
	Deepseek-V3	<b>57.93</b>	<b>29.27</b>	<b>33.54</b>	<b>35.37</b>	<b>39.03</b>	<b>66.46</b>	<b>36.59</b>	<b>37.80</b>	<b>37.80</b>	<b>44.66</b>
	Qwen-plus	0.00	3.05	4.27	4.88	3.05	19.51	7.93	5.49	7.93	10.22
ES	GPT-4o	34.92	31.54	31.37	29.88	31.93	39.02	30.02	31.88	29.69	32.65
	Deepseek-V3	<b>42.65</b>	<b>34.16</b>	<b>34.29</b>	<b>33.17</b>	<b>36.07</b>	<b>44.58</b>	<b>34.14</b>	<b>35.75</b>	<b>33.84</b>	<b>37.08</b>
	Qwen-plus	10.64	15.78	16.42	15.71	14.64	27.44	22.06	21.57	22.29	23.34

TABLE II: Comparisons between our decompilation LLMs with the baselines on HumanEval. \*: Our reproduction using LLM4Decompile’s original settings. †: Results are from the original paper. (%)

Metric	Model	HumanEval 32-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o +CIW	95.73	91.46	<b>92.07</b>	<b>93.90</b>	<b>93.29</b>	<b>97.56</b>	91.46	<b>94.51</b>	84.76	92.07
	Deepseek-V3 +CIW	<b>96.34</b>	<b>92.07</b>	89.63	90.85	92.22	<b>97.56</b>	87.80	85.37	78.66	87.35
	Qwen-plus +CIW	90.24	86.59	85.98	91.46	88.57	89.63	87.20	87.20	70.12	83.54
	LLM4Decompile 1.3B*	9.15	25.00	17.68	18.90	17.68	56.10	58.54	54.27	56.10	56.25
	LLM4Decompile 6.7B*	7.32	34.15	35.98	35.98	28.35	71.95	80.49	75.00	75.00	75.61
	FAE†	-	-	-	-	-	92.07	<b>93.29</b>	92.07	<b>93.90</b>	<b>92.84</b>
	CIM-1.3B +CIW	85.98	87.20	90.24	89.36	88.26	90.85	87.80	87.80	86.59	88.26
Re-exe	CIM-6.7B +CIW	89.02	90.85	90.24	92.07	90.55	93.29	91.46	93.29	92.07	92.53
	GPT-4o +CIW	32.32	28.05	24.39	24.39	27.29	45.12	29.27	23.17	19.51	29.27
	Deepseek-V3 +CIW	54.27	40.24	34.76	34.15	40.86	71.34	45.73	45.73	42.07	51.22
	Qwen-plus +CIW	22.56	17.07	15.85	17.68	18.29	26.83	13.41	12.20	12.80	16.31
	LLM4Decompile 1.3B*	0.00	1.22	0.61	0.61	0.61	25.61	10.37	7.32	9.76	13.26
	LLM4Decompile 6.7B*	0.61	1.22	0.61	0.61	0.76	39.02	26.22	30.49	28.05	30.94
	FAE†	-	-	-	-	-	71.95	53.66	48.78	45.73	55.03
ES	CIM-1.3B +CIW	65.24	34.15	35.37	31.71	41.62	71.34	39.63	42.07	40.24	48.32
	CIM-6.7B +CIW	<b>75.61</b>	<b>53.05</b>	<b>53.05</b>	<b>50.00</b>	<b>57.93</b>	<b>80.49</b>	<b>57.93</b>	<b>56.71</b>	<b>53.05</b>	<b>62.05</b>
	GPT-4o +CIW	40.81	35.47	36.28	35.41	36.99	44.48	36.45	36.86	34.11	37.98
	Deepseek-V3 +CIW	45.47	36.15	36.98	35.33	38.48	<b>51.36</b>	36.33	37.61	33.61	39.73
	Qwen-plus +CIW	39.71	34.80	35.76	34.37	36.16	43.59	35.75	35.14	31.63	36.53
	LLM4Decompile 1.3B*	14.05	13.42	11.23	11.08	12.45	30.14	19.89	18.81	20.26	22.27
	LLM4Decompile 6.7B*	29.03	20.52	21.80	21.15	23.12	41.42	32.07	32.03	31.68	34.30
	CIM-1.3B +CIW	47.45	35.88	37.03	36.33	39.17	47.56	35.67	37.64	37.31	39.55
	CIM-6.7B +CIW	<b>51.06</b>	<b>40.16</b>	<b>39.96</b>	<b>40.48</b>	<b>42.92</b>	49.16	<b>40.25</b>	<b>39.54</b>	<b>39.02</b>	<b>41.99</b>

bit benchmark.

### B. Main Results

The main experimental results can be seen at Table I-IV. Table I reports the decompilation results of general-purpose LLMs on the HumanEval without the assistance of **CIW**. Tables II and III present the decompilation results on the HumanEval and ExeBench: results for general-purpose LLMs are obtained with the assistance of **CIW**, while those for decompilation-specific models are either reproduced or cited from the original papers, depending on the availability of models and code. Unlike the evaluation metrics used in Tables

I, II and III, Table IV presents results based on the Pass@k metric. The results lead to the following key findings:

(1) By comparing Table I and Table II (i.e., vertical evaluation), we find that **CIW** significantly improves decompilation performance. On average, the three models achieve absolute gains of 15.26% on Rec-com, 9.16% on Re-exe, and 8.36% on ES. Notably, Qwen-plus shows the largest improvement on Rec-com (35.21%), GPT-4o on Re-exe (12.61%), and Qwen-plus again on ES (17.36%). These suggest that **CIW** contributes differently to various aspects of decompilation across different models.

(2) As shown in Tables II and III, in the horizontal eval-

TABLE III: Comparisons between our decompilation LLMs with the baselines on ExeBench. †: Results are from the original paper. (%)

Metric	Model	ExeBench 32-bit					ExeBench 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o +CIW	80.07	<b>85.25</b>	<b>84.80</b>	<b>85.13</b>	<b>83.81</b>	90.54	88.34	88.09	87.28	88.56
	Deepseek-V3 +CIW	84.30	81.30	80.15	81.41	81.79	<b>93.32</b>	<b>89.68</b>	<b>89.98</b>	<b>88.68</b>	<b>90.42</b>
	Qwen-plus +CIW	71.97	78.68	78.66	80.00	77.33	82.84	82.78	83.20	81.41	82.56
	CIM-1.3B +CIW	84.14	73.55	74.72	73.28	76.42	88.33	70.76	70.72	69.99	74.95
	CIM-6.7B +CIW	<b>86.87</b>	73.77	73.01	73.52	76.79	89.49	70.57	70.20	68.14	74.60
Re-exe	GPT-4o +CIW	30.87	30.20	28.27	29.08	29.61	43.99	25.61	23.61	22.06	28.82
	Deepseek-V3 +CIW	43.83	36.55	34.26	34.51	37.29	59.16	36.48	32.89	30.86	39.85
	Qwen-plus +CIW	22.05	22.54	21.22	21.28	21.77	32.96	19.85	17.33	16.33	21.62
	LLM4Decompile 1.3B†	-	-	-	-	-	17.86	13.62	13.20	13.28	14.49
	LLM4Decompile 6.7B†	-	-	-	-	-	22.89	16.60	16.18	16.25	17.98
	CIM-1.3B +CIW	56.09	35.33	33.94	32.49	39.46	65.20	36.32	33.34	32.55	41.85
	CIM-6.7B +CIW	<b>64.74</b>	<b>42.86</b>	<b>40.60</b>	<b>40.66</b>	<b>47.21</b>	<b>72.13</b>	<b>40.42</b>	<b>36.57</b>	<b>35.45</b>	<b>46.14</b>
ES	GPT-4o +CIW	52.02	40.53	38.68	38.44	42.42	46.53	41.24	39.85	40.16	41.95
	Deepseek-V3 +CIW	54.93	44.73	42.97	42.70	46.33	60.95	44.31	41.81	41.20	47.07
	Qwen-plus +CIW	43.62	40.39	39.15	39.06	40.56	50.98	39.87	38.59	37.70	41.79
	CIM-1.3B +CIW	66.14	50.02	49.03	48.03	53.31	67.58	50.27	49.21	48.28	53.84
	CIM-6.7B +CIW	<b>69.94</b>	<b>53.04</b>	<b>51.28</b>	<b>50.68</b>	<b>56.24</b>	<b>67.94</b>	<b>52.72</b>	<b>51.03</b>	<b>50.48</b>	<b>55.54</b>

TABLE IV: Comparisons between our decompilation LLMs with the baselines on HumanEval 64-bit. †: The results are from Nova [14] paper. (%)

Model	Pass @ 1					Pass @ 10				
	O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
GPT-4o†	21.34	18.29	14.48	13.05	16.79	29.94	26.74	21.42	19.88	24.50
LLM4Decompile 1.3B†	15.30	8.26	9.36	8.38	10.33	21.79	15.23	16.17	13.70	16.72
Nova 1.3B†	37.53	21.71	22.68	18.75	25.17	49.38	34.84	36.95	32.03	38.30
LLM4Decompile 6.7B†	29.97	19.05	20.46	18.32	21.95	40.40	27.75	28.85	28.51	31.38
Nova 6.7B†	48.78	30.58	30.85	27.23	34.36	57.47	47.45	43.03	39.68	46.91
CIM-1.3B +CIW	70.64	38.41	38.63	37.26	46.23	82.71	59.75	57.74	55.70	63.97
CIM-6.7B +CIW	<b>79.66</b>	<b>57.56</b>	<b>55.70</b>	<b>52.96</b>	<b>61.47</b>	<b>90.35</b>	<b>76.79</b>	<b>71.48</b>	<b>67.68</b>	<b>76.57</b>

uation, our 6.7B LLM achieves SOTA performance on Re-exe and the ES metric, with average improvements of 11.03% and 6.27% (on two datasets), respectively, over DeepSeek-V3 671B—a model with 100× more parameters. This highlights our model’s strength in preserving semantic consistency and functional correctness. While our Re-com score is slightly lower than that of some general-purpose LLMs, this is largely due to their strong contextual modeling capability, which enables them to produce syntactically complete code that may not faithfully reflect the original program logic.

(3) Table IV presents the results under the Pass@k evaluation metric. While Nova [14] currently achieves the best performance under the Pass@k metric, our model outperforms Nova at the same scale; notably, even our 1.3B model surpasses the performance of Nova’s 6.7B variant (i.e., our 46.23% vs Nova’s 34.36%, our 63.97% vs Nova’s 46.91%). Additionally, our Pass@1 results are closely aligned with the Re-exe scores reported in Table II (i.e., 46.23%↔48.32%, 61.47%↔62.05%). This is expected, as Re-exe evaluates the correctness of a single generation, while Pass@k further

reflects the stability and consistency of the model’s outputs across  $k$  attempts.

(4) From the results across all four tables, we observe a clear trend across all models: decompilation performance generally degrades as the compiler optimization level increases (O0→O1→O2→O3). For example, on the ExeBench 64-bit under the Re-exe metric, the average performance at O0 is 46.11%, O1 is 32.84%, O2 is 30.63%, and O3 is 28.68%. This decline is primarily due to the increasingly aggressive optimizations applied at higher levels (e.g., O2 and O3), such as expression rewriting, variable merging, dead code elimination, and control flow simplification. These transformations significantly alter the original function structure and reduce code readability, making it more challenging for models to accurately reconstruct the source code.

### C. Ablation Study

To evaluate the contribution of individual components in our proposed **CIW**, we conduct an ablation study (also referred to as vertical experiments). Specifically, for a given LLM,

TABLE V: Ablation study for CFG. (%)

Metric	Model	ExeBench 64-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	CIM-1.3B +CIW	<b>87.34</b>	<b>65.02</b>	<b>62.70</b>	<b>62.35</b>	<b>69.35</b>	91.77	88.54	89.17	87.90	89.35
	w/o cfg	87.02	63.14	61.44	61.18	68.19	90.51	90.45	90.45	89.17	90.15
	CIM-6.7B +CIW	87.23	62.75	59.77	61.81	67.89	<b>93.04</b>	<b>92.99</b>	<b>94.27</b>	<b>92.99</b>	<b>93.32</b>
	w/o cfg	85.98	62.06	60.29	60.57	67.23	92.41	91.08	88.54	90.45	90.62
Re-exe	CIM-1.3B +CIW	61.55	30.61	28.32	25.65	36.53	72.15	38.85	42.04	40.13	48.29
	w/o cfg	59.87	28.96	26.22	24.15	34.80	51.27	25.48	26.75	25.48	32.25
	CIM-6.7B +CIW	<b>67.72</b>	<b>32.33</b>	<b>28.62</b>	<b>27.95</b>	<b>39.04</b>	<b>80.38</b>	<b>58.60</b>	<b>57.32</b>	<b>53.50</b>	<b>62.45</b>
	w/o cfg	58.61	28.75	25.89	24.15	34.35	65.19	34.39	41.40	37.58	44.64
ES	CIM-1.3B +CIW	67.66	50.76	49.29	48.66	54.09	48.05	36.06	38.27	37.92	40.08
	w/o cfg	66.83	50.60	48.53	47.61	53.39	45.28	33.38	31.70	31.14	35.38
	CIM-6.7B +CIW	<b>68.65</b>	<b>52.30</b>	<b>50.95</b>	<b>49.85</b>	<b>55.44</b>	<b>49.86</b>	<b>40.97</b>	<b>40.09</b>	<b>39.56</b>	<b>42.62</b>
	w/o cfg	66.69	50.00	48.40	47.44	53.13	45.59	34.83	34.72	34.88	37.51

TABLE VI: Ablation study for data mapping table. (%)

Model	Metric	ExeBench 64-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	CIM-1.3B +CIW	84.37	60.67	58.89	57.10	65.26	91.77	88.54	89.17	87.90	89.35
	w/o DM	80.18	60.39	58.33	58.45	64.34	<b>95.00</b>	90.24	92.31	90.57	92.03
	CIM-6.7B +CIW	<b>86.86</b>	<b>62.64</b>	<b>60.56</b>	59.25	<b>67.33</b>	87.50	<b>92.68</b>	<b>94.87</b>	<b>94.34</b>	<b>92.35</b>
	w/o DM	84.03	56.74	59.44	<b>61.39</b>	65.40	<b>95.00</b>	90.24	84.62	84.91	88.69
Re-exe	CIM-1.3B +CIW	55.83	23.31	14.72	13.40	26.82	72.15	38.85	42.04	40.13	48.29
	w/o DM	52.89	21.35	14.44	12.60	25.32	55.00	26.83	33.33	32.08	36.81
	CIM-6.7B +CIW	<b>65.69</b>	<b>28.93</b>	<b>18.89</b>	<b>17.96</b>	<b>32.87</b>	<b>75.00</b>	<b>48.78</b>	<b>43.59</b>	<b>43.40</b>	<b>52.69</b>
	w/o DM	60.93	22.47	17.50	17.43	29.58	55.00	41.46	41.59	41.51	44.89
ES	CIM-1.3B +CIW	68.43	47.89	39.43	38.28	48.51	42.94	33.10	<b>36.59</b>	36.32	37.24
	w/o DM	66.48	47.63	38.80	38.45	47.84	45.14	34.58	33.12	35.65	37.12
	CIM-6.7B +CIW	<b>71.61</b>	<b>51.81</b>	42.98	<b>42.91</b>	<b>52.33</b>	47.41	36.40	35.78	37.24	39.21
	w/o DM	70.38	48.24	<b>44.18</b>	42.28	51.27	<b>48.16</b>	<b>37.76</b>	35.02	<b>37.51</b>	<b>39.61</b>

we selectively remove either the CFG or the data mapping table (DM) from the input and observe the resulting impact on decompilation performance. The Tables V and VI present the results<sup>2</sup> of the ablation study on our LLMs. More ablations results are shown in Tables VII, IX, VIII, X.

From the results, we observe two key findings: (1) Both CFG and DM can improve model performance, but in complementary ways. Removing CFG leads to consistent performance drops across all metrics—especially on execution-related scores—highlighting the importance of structural control-flow information in preserving program semantics. In contrast, removing DM causes smaller but still noticeable degradations, particularly in Re-exe and ES metrics, suggesting that data-level correspondences enhance variable-level alignment and execution fidelity. (2) Models that integrate both CFG and DM

demonstrate more robust performance across different compiler optimization levels (O0–O3). In contrast, models with either component removed show greater performance variance and occasional regressions. This indicates that combining structural (CFG) and symbolic (DM) inductive biases enables better generalization to real-world, compiler-transformed binaries.

#### D. Case Study

To qualitatively assess CIM’s performance, we conduct the case study. Full case study results and visual comparisons are provided in follows.

##### (1) Effectiveness of CIM

We evaluate CIM’s effectiveness through a qualitative comparison with three baselines: GPT-4o with CIW, DeepSeek-V3 with CIW, and Hex-Rays, using the -O1-compiled 32-bit binary of function *r\_deriv* from test set of ExeBench (Fig. 4). While all methods correctly recover data object (which have

<sup>2</sup>We conduct the corresponding ablation experiments only on samples where the number of CFG blocks is greater than one or where data sections are present.

```

1 __attribute__((used)) static int r_deriv(struct
2     SN_env * z) {
3     int among_var;
4     z->ket = z->c;
5     among_var = find_among_b(z, a_2, 25);
6     if (!among_var) return 0;
7     z->bra = z->c;
8     switch (among_var) {
9         case 1:
10            { int ret = r_R2(z);
11              if (ret <= 0) return ret;
12            }
13            { int ret = slice_del(z);
14              if (ret < 0) return ret;
15            }
16            break;
17        case 2:
18            { int ret = slice_from_s(z, 3, s_9)
19              ;
20              if (ret < 0) return ret;
21            }
22            break;
23        case 3:
24            { int ret = slice_from_s(z, 3, s_10)
25              ;
26              if (ret < 0) return ret;
27            }
28            break;
29        case 4:
30            { int ret = slice_from_s(z, 4, s_11)
31              ;
32              if (ret < 0) return ret;
33            }
34            break;
35        case 5:
36            { int ret = slice_from_s(z, 5, s_12)
37              ;
38              if (ret < 0) return ret;
39            }
40            break;
41        case 6:
42            { int ret = slice_from_s(z, 3, s_13)
43              ;
44              if (ret < 0) return ret;
45            }
46            break;
47        }
48    }
49    return 1;
50 }

```

(a) Source code

```

1 int r_en_ending_00000001(struct
2     SN_env * z) {
3     int among_var;
4     z->ket = z->c;
5     among_var = find_among_b(z, a_2,
6     25);
7     if (!among_var) return 0;
8     z->bra = z->c;
9     switch (among_var) {
10        case 0: return 0;
11        case 1:
12            { int ret = r_R2(z);
13              if (ret <= 0) return
14              ret;
15            }
16            { int ret = slice_del(z);
17              if (ret < 0) return ret
18              ;
19            }
20            break;
21        case 2:
22            { int ret = slice_from_s(
23              z, 3, s_9);
24              if (ret < 0) return ret
25              ;
26            }
27            break;
28        case 3:
29            { int ret = slice_from_s(
30              z, 3, s_10);
31              if (ret < 0) return ret
32              ;
33            }
34            break;
35        case 4:
36            { int ret = slice_from_s(
37              z, 4, s_11);
38              if (ret < 0) return ret
39              ;
40            }
41            break;
42        case 5:
43            { int ret = slice_from_s(
44              z, 5, s_12);
45              if (ret < 0) return ret
46              ;
47            }
48            break;
49        case 6:
50            { int ret = slice_from_s(
51              z, 3, s_13);
52              if (ret < 0) return ret
53              ;
54            }
55            break;
56        }
57    }
58    return 1;
59 }

```

(b) CIM-6.7B + CIW

```

1 __attribute__((fastcall)) function
2     (int z) {
3     int* esi = (int*)z;
4     int ebp = *esi;
5     esi[2] = ebp;
6     int among_var = find_among_b(es,
7     (const char*)(a_2 - 3BE68h),
8     25);
9     int ebx = among_var;
10    if (among_var) {
11        esi[1] = ebp;
12        switch (among_var) {
13            case 1:
14                if (r_R2(es) > 0) {
15                    slice_del(es);
16                } else {
17                    among_var = 0;
18                }
19                break;
20            case 2:
21                among_var =
22                slice_from_s(es,
23                3, (const char*)(
24                s_9 - 3BE68h));
25                break;
26            case 3:
27                among_var =
28                slice_from_s(es,
29                3, (const char*)(
30                s_10 - 3BE68h));
31                break;
32            case 4:
33                among_var =
34                slice_from_s(es,
35                4, (const char*)(
36                s_11 - 3BE68h));
37                break;
38            case 5:
39                among_var =
40                slice_from_s(es,
41                5, (const char*)(
42                s_12 - 3BE68h));
43                break;
44            case 6:
45                among_var =
46                slice_from_s(es,
47                3, (const char*)(
48                s_13 - 3BE68h));
49                break;
50            default:
51                among_var = 1;
52                break;
53        }
54    }
55    return ebx;
56 }

```

(c) DeepSeek-V3 + CIW

```

1 int function_name(int* z) {
2     int ebp = z[0];
3     z[2] = ebp;
4     int e = 25;
5     const char* d = "a_2";
6     int ebx = find_among_b(z, d, e);
7     int among_var = 0;
8     if (ebx != 0) {
9         z[1] = ebp;
10        if (among_var <= 6) {
11            switch (among_var) {
12                case 1:
13                    if (r_R2(z) <= 0) {
14                        among_var = 0;
15                    } else {
16                        if (slice_del(z) != 0)
17                            among_var = 0;
18                    }
19                    break;
20                case 2:
21                    if (slice_from_s(z, 3, "s_9
22                    ") != 0) {
23                        among_var = 0;
24                    }
25                    break;
26                case 3:
27                    if (slice_from_s(z, 3, "
28                    s_10") != 0) {
29                        among_var = 0;
30                    }
31                    break;
32                case 4:
33                    if (slice_from_s(z, 4, "
34                    s_11") != 0) {
35                        among_var = 0;
36                    }
37                    break;
38                case 5:
39                    if (slice_from_s(z, 5, "
40                    s_12") != 0) {
41                        among_var = 0;
42                    }
43                    break;
44                case 6:
45                    if (slice_from_s(z, 3, "
46                    s_13") != 0) {
47                        among_var = 0;
48                    }
49                    break;
50                default:
51                    among_var = 1;
52                    break;
53            }
54        } else {
55            among_var = 1;
56        }
57    }
58    return ebx;
59 }

```

(d) GPT-4o + CIW

```

1 int r_deriv(SN_env *z)
2 {
3     int v1;
4     int v2;
5     int v3;
6     int v4;
7     int v5;
8
9     v1 = z->c;
10    z->ket = z->c;
11    v2 = find_among_b(z, a_2, 25);
12    v3 = v2;
13    if (v2)
14    {
15        z->bra = v1;
16        switch (v2)
17        {
18            case 1:
19                v4 = r_R2(z);
20                if (v4 <= 0)
21                {
22                    v3 = v4;
23                }
24                else
25                {
26                    v5 = slice_del(z);
27                    if (v5 < 0)
28                        v3 = v5;
29                }
30                break;
31            case 2:
32                v3 = slice_from_s(z, 3, s_9);
33                if (v3 >= 0)
34                    v3 = 1;
35                break;
36            case 3:
37                v3 = slice_from_s(z, 3, s_10)
38                ;
39                if (v3 >= 0)
40                    v3 = 1;
41                break;
42            case 4:
43                v3 = slice_from_s(z, 4, s_11)
44                ;
45                if (v3 >= 0)
46                    v3 = 1;
47                break;
48            case 5:
49                v3 = slice_from_s(z, 5, s_12)
50                ;
51                if (v3 >= 0)
52                    v3 = 1;
53                break;
54            case 6:
55                v3 = slice_from_s(z, 3, s_13)
56                ;
57                if (v3 >= 0)
58                    v3 = 1;
59                break;
60            default:
61                v3 = 1;
62                break;
63        }
64    }
65    return v3;
66 }

```

(e) Hex-Rays

Fig. 4: An example of readable decompilation. Presented are the source code (a) of the case sample alongside the decompilations of CIM-6.7B with CIW (b), DeepSeek-V3 with CIW (c), GPT-4o with CIW (d), and Hex-Rays (e).

---

```

1 void runHitch() {
2     float lr7;
3
4     lr7 = vexRT[JOY_BTN_HITCH_UP] - vexRT[
        JOY_BTN_HITCH_DOWN];
5     if(lr7 != 0.0)
6         motor[servoHitch] = lr7 * 128;
7 }

```

---

(a) Source code

---

```

1 void moveHitch() {
2     float hitchSpeed = vexRT[JOY_BTN_HITCH_UP
        ] - vexRT[JOY_BTN_HITCH_DOWN];
3
4     if(hitchSpeed != 0.0)
5         motor[servoHitch] = hitchSpeed * 128;
6     else
7         motor[servoHitch] = 0;
8 }

```

---

(c) CIM-6.7B + CIW w/o CFG

---

```

1 void moveHitch() {
2     float hitchPower = vexRT[JOY_BTN_HITCH_UP
        ] - vexRT[JOY_BTN_HITCH_DOWN];
3
4     if(hitchPower != 0.0) {
5         motor[servoHitch] = hitchPower * 128;
6     }
7 }

```

---

(b) CIM-6.7B + CIW

---

```

1 void moveHitch() {
2     float hitchPower = vexRT[JOY_BTN_HITCH_UP
        ] - vexRT[JOY_BTN_HITCH_DOWN];
3
4     if(hitchPower != 0.0) {
5         motor[servoHitch] = hitchPower * 1.0;
6     }
7 }

```

---

(d) CIM-6.7B + CIW w/o DM

Fig. 5: Qualitative example of the decompilations made by CIM with CIW and its variations that ablate different components. Presented are the source code (a) of the case sample alongside the decompilations of CIM-6.7B + CIW (b), CIM-6.7B + CIW w/o CFG (c), and CIM-6.7B + CIW w/o DM (d).

no data mapping between memory addresses and data segments), CIM with CIW uniquely achieves both correctness and enhanced readability in reconstructing program logic.

The baseline methods fail to correctly reconstruct program logic during decompilation. As shown in Fig. 4(d), GPT-4o with CIW makes three critical errors: (1) hardcoding the `switch` condition (line 10) to 0 instead of deriving it from function `find_among_b`’s return value, (2) incorrectly checking for non-zero rather than negative values from function `slice_from_s` in cases 2–6 of `switch`, and (3) misassociating function `r_deriv`’s return logic solely with function `find_among_b` (lines 6,54). DeepSeek-V3 with CIW exhibits similar issues, failing to reconstruct both the `switch` logic and proper return behavior (lines 6-7,41 in Fig. 4(c)). Even Hex-Rays(Fig. 4(e)), while correctly handling return logic, inverts the condition in cases 2-6 of `switch` to check for the return value of function `slice_from_s`  $\geq 0$  instead of  $<0$ . These consistent failures across diverse methods highlight the difficulty of accurately recovering program logic from low-level assembly code.

In contrast, CIM with CIW (Fig. 4(b)) successfully addresses these failures through fine-tuning on a novel dataset (CID). CIM-6.7B with CIW achieves fully accurate reconstruction of: (1) the `switch` entry condition, (2) case logic in branches 2-6, and (3) function `r_deriv` return logic. Beyond basic correctness, CIM demonstrates superior code quality through: (a) adoption of modern coding conventions like "Early Return" (line 5) to reduce control flow complexity, (b) generation of meaningful variable names and types, and (c) elimination of unnecessary procedural variables. This dual achievement of logical accuracy and code quality stems from CIM’s core innovation - moving beyond simple instruction translation to explicitly model programmer intent through

comprehensive control flow analysis during training.

## (2) Reasonableness of CIW

Our qualitative evaluation validates the CIW’s design by analyzing CIM-6.7B’s performance on the function `runHitch` (from ExeBench, compiled with -O0 for 32-bit architecture) through systematic ablation studies (Figure 5). The results reveal that CIM-6.7B with CIW produces correct decompilations, models with ablated variants exhibit distinct failure modes that underscore CIW’s careful component integration. When CIW without CFG information, CIM-6.7B introduces structural flaws, including an erroneous `else` branch that corrupts the value of `motor[servoHitch]` (Fig. 5(c), lines 7-8), demonstrating that program logic reconstruction fundamentally requires CFG guidance. Conversely, the CIW without data mapping tables version maintains correct program logic but fails at data recovery, notably converting constant 128 to 1.0 (Fig. 5(d), line 5), revealing that data mapping tables is critical for data recovery.

These results validate CIW’s core design principle: high-quality decompilation demands simultaneous attention to both macro-level structural accuracy (via CFG analysis) and micro-level value precision (through data recovery). The success of CIM-6.7B with complete CIW, contrasted with the failures of ablated variants, demonstrates the necessity of its integrated component of CIW for producing behaviorally accurate decompilations.

## V. CONCLUSION

Despite recent advances, end-to-end decompilation with LLMs still suffers from limited control flow awareness, weak data context modeling, and high computational costs. This

paper addresses these challenges by introducing a structure-aware and data-informed decompilation scheme that integrates CFGs and data mapping tables into LLM inputs. We build a dedicated dataset and propose CIM, a lightweight domain-specific LLM, and demonstrate that it achieves state-of-the-art performance in re-executability, Pass@k, and readability, while significantly reducing computational costs to enable efficient and locally deployable inference.

## VI. FUTURE WORK

While this work explores the integration of control flow graphs (CFGs) and data mapping tables into large language models for decompilation, several important challenges remain open for future investigation.

First, the potential of data flow graphs (DFGs) in enhancing decompilation quality has not been fully explored. Incorporating DFGs could provide complementary semantic and dependency information that may further improve the model's understanding of binary code.

Second, our experiments indicate a significant drop in decompilation performance on binaries compiled with higher optimization levels. Understanding the underlying reasons behind this degradation and designing models that are more robust to aggressive compiler optimizations are promising directions for future research.

## REFERENCES

- [1] C. Kruegel, W. Robertson, and et al., "Detecting kernel-level rootkits through binary analysis," in *Proceedings of the 20th Annual Computer Security Applications Conference*, ser. ACSAC '04. USA: IEEE Computer Society, 2004, p. 91–100. [Online]. Available: <https://doi.org/10.1109/CSAC.2004.19>
- [2] C. Fu, H. Chen, and et al., "Coda: an end-to-end neural program decompiler," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [3] Y. David, U. Alon, and et al., "Neural reverse engineering of stripped binaries using augmented control flow graphs," *Proc. ACM Program. Lang.*, vol. 4, no. OOPSLA, Nov. 2020. [Online]. Available: <https://doi.org/10.1145/3428293>
- [4] M. Egele, T. Scholte, and et al., "A survey on automated dynamic malware-analysis techniques and tools," *ACM Comput. Surv.*, vol. 44, no. 2, Mar. 2008. [Online]. Available: <https://doi.org/10.1145/2089125.2089126>
- [5] L. Szekeres, M. Payer, and et al., "Sok: Eternal war in memory," in *2013 IEEE Symposium on Security and Privacy*, 2013, pp. 48–62.
- [6] Z. J. Zhou, R. C. Dong, J. H. Jiang, and W. H. Zhang, "Survey on binary code security techniques," *Computer Systems and Applications*, vol. 32, no. 1, pp. 1–11, 2023, in Chinese. [Online]. Available: <http://www.c-s-a.org.cn/1003-3254/8848.html>
- [7] G. Bossert, F. Guihéry, and et al., "Towards automated protocol reverse engineering using semantic information," in *Proceedings of the 9th ACM Symposium on Information, Computer and Communications Security*, ser. ASIA CCS '14, New York, NY, USA, 2014, p. 51–62. [Online]. Available: <https://doi.org/10.1145/2590296.2590346>
- [8] Hex-Rays, "Hex-rays decompiler," 2025, accessed: 2025-01-14. [Online]. Available: <https://hex-rays.com/decompiler/>
- [9] Ghidra, "Ghidra software reverse engineering framework," 2025, accessed: 2025-01-14. [Online]. Available: <https://ghidra-sre.org/>
- [10] Y. Cao, R. Liang, and et al., "Boosting neural networks to decompile optimized binaries," in *Proceedings of the 38th Annual Computer Security Applications Conference*, ser. ACSAC '22, New York, NY, USA, 2022, p. 508–518. [Online]. Available: <https://doi.org/10.1145/3564625.3567998>
- [11] C. Zhu, Z. Li, and et al., "Tygr: Type inference on stripped binaries using graph neural networks," in *USENIX Security Symposium*, 2024.
- [12] R. Liang, Y. Cao, and et al., "Neutron: an attention-based neural decompiler," *Cybersecurity*, vol. 4, p. 5, 03 2021.
- [13] H. Tan, Q. Luo, and et al., "Llm4decompile: Decompiling binary code with large language models," 2024.
- [14] N. Jiang, C. Wang, and et al., "Nova: Generative language models for assembly code with hierarchical attention and contrastive learning," in *Submitted to The Thirteenth International Conference on Learning Representations*, 2024, under review.
- [15] Y. Feng, D. Teng, and et al., "Self-constructed context decompilation with fined-grained alignment enhancement," 2024.
- [16] W. K. Wong, H. Wang, and et al., "Refining decompiled c code with large language models," 2023.
- [17] P. Hu, R. Liang, and et al., "Degpt: Optimizing decompiler output with llm," in *Proceedings 2024 Network and Distributed System Security Symposium (2024)*. <https://api.semanticscholar.org/CorpusID.267622140>, 2024.
- [18] X. She, Y. Zhao, and et al., "Wadec: Decompiling webassembly using large language model," in *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '24, New York, NY, USA, 2024, p. 481–492. [Online]. Available: <https://doi.org/10.1145/3691620.3695020>
- [19] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar et al., "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [20] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim et al., "Starcode: may the source be with you!" *arXiv preprint arXiv:2305.06161*, 2023.
- [21] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever et al., "Improving language understanding by generative pre-training," 2018.
- [22] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman et al., "Evaluating large language models trained on code," *arXiv preprint arXiv:2107.03374*, 2021.
- [23] Z. Yu, W. Zheng, and et al., "(supplementary) codecmr: Cross-modal retrieval for function-level binary source code matching," 2020.
- [24] Z. Yu, R. Cao, and et al., "Order matters: Semantic-aware neural networks for binary code similarity detection," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, pp. 1145–1152, Apr. 2020.
- [25] H. Gao, S. Cheng, and et al., "A lightweight framework for function name reassignment based on large-scale stripped binaries," in *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2021, New York, NY, USA, 2021, p. 607–619. [Online]. Available: <https://doi.org/10.1145/3460319.3464804>
- [26] J. Armengol-Estapé, J. Woodruff, and et al., "Exebench: an ml-scale dataset of executable c functions," in *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, ser. MAPS 2022. New York, NY, USA: Association for Computing Machinery, 2022, p. 50–59. [Online]. Available: <https://doi.org/10.1145/3520312.3534867>
- [27] S. Hex-Rays, "Ida pro: a cross-platform multi-processor disassembler and debugger," 2014.
- [28] Binutils, "objdump," 2025, [Online]. [Online]. Available: <https://www.gnu.org/software/binutils/>
- [29] Y. Feng, B. Li, X. Shi, Q. Zhu, and W. Che, "Ref decompile: Relabeling and function call enhanced decompile," *arXiv preprint arXiv:2502.12221*, 2025.
- [30] M. Chen, J. Tworek, and et al., "Evaluating large language models trained on code," 2021. [Online]. Available: <https://arxiv.org/abs/2107.03374>
- [31] J. Armengol-Estapé, J. Woodruff, and et al., "Slade: A portable small language model decompiler for optimized assembly," in *Proceedings of the 2024 IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO '24. IEEE Press, 2024, p. 67–80. [Online]. Available: <https://doi.org/10.1109/CGO57630.2024.10444788>
- [32] X. Du, M. Liu, K. Wang, H. Wang, J. Liu, Y. Chen, J. Feng, C. Sha, X. Peng, and Y. Lou, "Evaluating large language models in class-level code generation," in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, ser. ICSE '24. New York, NY, USA: Association for Computing Machinery, 2024.
- [33] D. S. Katz, J. Ruchtli, and et al., "Using recurrent neural networks for decompilation," in *2018 IEEE 25th International Conference on*

- [34] I. Hosseini and B. Dolan-Gavitt, “Beyond the c: Retargetable decompilation using neural machine translation,” in *Proceedings 2022 Workshop on Binary Analysis Research*, ser. BAR 2022. Internet Society, 2022.
- [35] OpenAI, J. Achiam, S. Adler, and et al., “Gpt-4 technical report,” 2024. [Online]. Available: <https://arxiv.org/abs/2303.08774>
- [36] DeepSeek-AI, A. Liu, B. Feng, and et al., “Deepseek-v3 technical report,” *ArXiv*, vol. abs/2412.19437, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:275118643>
- [37] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang *et al.*, “Qwen technical report,” *arXiv preprint arXiv:2309.16609*, 2023.
- [38] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, “Lora: Low-rank adaptation of large language models,” *ICLR*, vol. 1, no. 2, p. 3, 2022.
- [39] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica, “Efficient memory management for large language model serving with pagedattention,” in *Proceedings of the 29th Symposium on Operating Systems Principles*, ser. SOSP '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 611–626. [Online]. Available: <https://doi.org/10.1145/3600006.3613165>

## VII. ETHICS CONSIDERATIONS

Our work focuses on leveraging large language models to improve the quality and efficiency of software decompilation, with the aim of supporting software maintenance, legacy system modernization, and software security auditing. We believe this research has positive societal impacts, such as enabling better understanding of binaries lacking source code, assisting vulnerability discovery, and promoting long-term software sustainability in critical systems.

However, we also acknowledge potential negative impacts. Decompilation technologies, especially when enhanced by powerful machine learning models, could be misused for reverse engineering proprietary software, circumventing software protections, or developing malware. To mitigate such risks, we explicitly position our research within the context of security auditing and legitimate software analysis, and strongly discourage malicious applications.

Regarding data and model ethics, all training data used in this work were sourced from public, permissively licensed datasets (e.g., open-source software with appropriate licenses). We performed due diligence to ensure that the data do not contain sensitive or proprietary content. Our model is trained solely on data intended for research and educational use.

While we recognize the value of open-source practices in promoting reproducibility and community collaboration, we are also mindful of the potential for misuse. If the model and code are released, we will consider applying responsible open-sourcing strategies, such as access limitations, usage licenses, or model cards outlining appropriate use cases and risks.

To discourage unethical use, we advocate for community norms that promote transparency, security-focused research, and legal compliance. We encourage future researchers and practitioners building upon our work to follow responsible disclosure practices and to engage with legal and ethical review processes when appropriate.

TABLE VII: Ablation study on HumanEval benchmark for CFG. (%)

Metric	Model	HumanEval 32-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o+CIW	95.57	91.08	91.72	93.63	93.00	97.47	91.08	94.27	84.71	91.88
	w/o cfg	93.04	93.63	93.63	92.63	93.17	96.20	92.36	92.36	87.90	92.21
	Deepseek-V3+CIW	97.47	91.72	89.81	84.71	90.93	94.94	88.54	85.35	82.17	87.75
	w/o cfg	93.67	92.99	91.09	90.45	92.05	95.57	89.81	87.90	83.44	89.18
	Qwen-plus+CIW	90.51	86.62	86.62	91.08	88.71	89.24	86.62	87.26	69.43	83.14
	w/o cfg	87.34	87.26	84.08	89.17	86.96	87.97	83.44	86.62	73.25	82.82
	CIM-1.3B+CIW	86.71	87.90	91.08	90.45	89.04	91.77	88.54	89.17	87.90	89.35
	w/o cfg	90.51	85.35	86.62	87.26	87.44	90.51	90.45	90.45	89.17	90.15
	CIM-6.7B+CIW	89.87	91.72	90.45	92.36	91.10	93.04	92.99	94.27	92.99	93.32
	w/o cfg	89.24	88.54	85.35	86.62	87.44	92.41	91.08	88.54	90.45	90.62
Re-exe	GPT-4o+CIW	31.01	26.11	23.57	22.29	25.75	43.67	27.39	20.38	17.20	27.16
	w/o cfg	24.68	15.29	14.01	12.74	16.68	37.97	12.74	13.38	9.55	18.41
	Deepseek-V3+CIW	50.63	33.76	29.94	31.85	36.55	68.99	42.68	43.31	38.22	48.30
	w/o cfg	50.00	23.57	28.66	29.30	32.88	67.62	33.76	36.94	30.57	42.22
	Qwen-plus+CIW	20.89	14.65	14.01	15.29	16.21	25.32	10.83	10.19	10.83	14.29
	w/o cfg	15.19	8.28	9.55	6.37	9.85	21.52	7.64	8.92	4.46	10.64
	CIM-1.3B+CIW	65.19	33.12	35.03	31.21	41.14	72.15	38.85	42.04	40.13	48.29
	w/o cfg	50.63	26.11	20.38	22.29	29.85	51.27	25.48	26.75	25.48	32.25
	CIM-6.7B+CIW	75.32	52.87	52.87	49.68	57.69	80.38	58.60	57.32	53.50	62.45
	w/o cfg	65.82	35.03	37.58	38.22	44.16	65.19	34.39	41.40	37.58	44.64
ES	GPT-4o+CIW	41.51	36.02	36.74	35.84	37.53	45.07	36.93	37.37	34.53	38.48
	w/o cfg	39.06	35.64	35.40	36.04	36.54	45.08	36.21	36.52	35.48	38.32
	Deepseek-V3+CIW	46.37	36.69	37.62	35.84	39.13	52.33	36.62	38.04	34.11	40.28
	w/o cfg	43.98	36.12	37.77	36.99	38.72	50.99	38.32	39.80	36.54	41.41
	Qwen-plus+CIW	40.23	35.15	36.20	34.69	36.57	44.09	36.10	35.43	31.66	36.82
	w/o cfg	39.17	35.02	35.00	34.02	35.80	42.50	34.15	34.43	32.07	35.79
	CIM-1.3B+CIW	48.19	36.44	37.68	36.94	39.81	48.05	36.06	38.27	37.92	40.08
	w/o cfg	44.07	30.66	30.63	31.26	34.16	45.28	33.38	31.70	31.14	35.38
	CIM-6.7B+CIW	51.60	40.73	40.57	41.12	43.51	49.86	40.97	40.09	39.56	42.62
	w/o cfg	46.57	36.87	35.17	33.94	38.14	45.59	34.83	34.72	34.88	37.51

TABLE VIII: Ablation study on ExeBench benchmark for CFG. (%)

Metric	Model	ExeBench 32-bit					ExeBench 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o+CIW	77.98	79.94	80.17	80.80	79.72	87.27	83.63	82.72	91.91	83.88
	w/o cfg	76.54	76.90	76.58	77.15	76.79	84.81	82.00	80.99	78.86	81.67
	Deepseek-V3+CIW	80.82	76.88	76.27	76.44	77.60	90.63	84.79	85.71	84.47	86.40
	w/o cfg	78.94	75.45	77.22	78.43	77.51	89.48	76.61	75.28	72.05	78.36
	Qwen-plus+CIW	68.24	71.11	70.19	73.01	70.64	77.90	75.45	76.46	75.44	76.31
	w/o cfg	68.14	67.01	66.33	66.12	66.90	78.86	67.44	67.50	65.31	69.53
	CIM-1.3B+CIW	83.60	65.46	64.94	64.83	69.71	87.34	65.02	62.70	62.35	69.35
	w/o cfg	80.86	67.60	69.15	66.05	70.91	87.02	63.14	61.44	61.18	68.19
	CIM-6.7B+CIW	85.50	67.20	65.43	65.31	70.86	87.23	62.75	59.77	61.81	67.89
	w/o cfg	84.00	63.39	63.19	63.55	68.53	85.98	62.06	60.29	60.57	67.23
Re-exe	GPT-4o+CIW	24.20	20.91	17.65	19.20	20.49	35.17	18.29	14.78	14.23	20.62
	w/o cfg	22.19	17.44	14.76	14.91	17.33	30.88	16.12	14.09	13.16	18.56
	Deepseek-V3+CIW	38.25	27.92	25.05	24.63	28.96	50.07	29.79	24.31	23.63	31.95
	w/o cfg	32.94	24.28	22.46	22.03	25.43	46.13	22.45	18.64	17.02	26.06
	Qwen-plus+CIW	17.23	13.32	12.30	12.48	13.83	25.50	12.75	10.02	10.37	14.66
	w/o cfg	15.97	11.70	10.03	10.16	11.97	23.49	8.60	7.15	6.64	11.47
	CIM-1.3B+CIW	53.58	29.21	27.54	26.82	34.29	61.55	30.61	28.32	25.65	36.53
	w/o cfg	50.79	27.56	27.50	24.30	32.54	59.87	28.96	26.22	24.15	34.80
	CIM-6.7B+CIW	61.90	33.91	30.31	29.82	38.98	67.72	32.33	28.62	27.95	39.04
	w/o cfg	53.87	28.11	25.89	25.25	33.28	58.61	28.75	25.89	24.15	34.35
ES	GPT-4o+CIW	44.17	36.34	35.10	34.74	37.59	48.03	35.27	33.16	33.30	37.44
	w/o cfg	44.01	35.33	34.09	34.05	36.87	48.48	34.82	33.49	34.70	37.87
	Deepseek-V3+CIW	52.13	38.92	37.05	36.47	41.14	55.97	37.92	35.20	34.78	40.97
	w/o cfg	52.04	38.51	36.80	36.45	40.95	55.91	37.02	35.30	33.66	40.47
	Qwen-plus+CIW	42.05	35.31	33.87	33.86	36.27	47.49	34.45	33.17	31.73	36.71
	w/o cfg	40.65	34.31	33.62	33.00	35.40	47.44	33.61	32.35	31.41	36.20
	CIM-1.3B+CIW	65.56	48.04	46.85	45.19	51.41	67.66	50.76	49.29	48.66	54.09
	w/o cfg	65.20	48.87	46.83	45.87	51.69	66.83	50.60	48.53	47.61	53.39
	CIM-6.7B+CIW	68.19	51.02	48.64	48.16	54.00	68.65	52.30	50.95	49.85	55.44
	w/o cfg	66.37	50.07	48.05	47.39	52.97	66.69	50.00	48.40	47.44	53.13

TABLE IX: Ablation study on HumanEval benchmark for Data Mapping(DM). (%)

Metric	Model	HumanEval 32-bit					HumanEval 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o+CIW	96.77	89.47	91.67	94.44	93.09	90.00	87.80	84.62	73.58	84.00
	w/o DM	93.55	84.21	83.33	63.89	81.25	85.00	80.49	69.23	67.92	75.66
	Deepseek-V3+CIW	100.00	92.11	88.89	83.33	91.08	92.50	75.61	74.36	71.70	78.54
	w/o DM	100.00	89.47	80.56	83.33	88.34	100.00	87.80	79.49	73.58	85.22
	Qwen-plus+CIW	80.65	86.84	83.33	83.30	83.53	80.00	78.05	66.67	47.17	67.97
	w/o DM	83.87	71.05	72.22	77.78	76.23	75.00	60.98	56.41	58.49	62.72
	CIM-1.3B+CIW	90.32	89.47	97.22	91.67	92.17	91.77	88.54	89.17	87.90	89.35
	w/o DM	90.32	89.47	88.89	94.44	90.78	95.00	90.24	92.31	90.57	92.03
	CIM-6.7B+CIW	96.77	97.32	94.44	97.22	96.44	87.50	92.68	94.87	94.34	92.35
	w/o DM	90.32	94.74	94.44	91.67	92.79	95.00	90.24	84.62	84.91	88.69
Re-exe	GPT-4o+CIW	16.13	15.79	11.11	8.33	12.84	22.50	19.51	23.08	13.21	19.58
	w/o DM	3.23	0.00	2.78	2.78	2.20	10.00	9.76	5.13	7.55	8.11
	Deepseek-V3+CIW	35.48	13.16	27.78	22.22	24.66	40.00	34.15	33.33	32.08	34.89
	w/o DM	16.13	2.63	8.33	0.00	6.77	50.00	17.07	30.77	22.64	30.12
	Qwen-plus+CIW	19.35	7.89	5.56	13.89	11.67	10.00	12.20	5.13	7.55	8.72
	w/o DM	9.68	5.26	2.78	2.78	5.13	7.50	0.00	0.00	0.00	1.88
	CIM-1.3B+CIW	70.97	18.42	22.22	25.00	34.15	72.15	38.85	42.04	40.13	48.29
	w/o DM	45.16	23.68	27.78	33.33	32.49	55.00	26.83	33.33	32.08	36.81
	CIM-6.7B+CIW	77.42	28.95	38.89	44.44	47.43	75.00	48.78	43.59	43.40	52.69
	w/o DM	51.61	23.68	30.56	30.56	34.10	55.00	41.46	41.59	41.51	44.89
ES	GPT-4o+CIW	42.59	36.85	37.28	34.98	37.93	41.59	36.15	37.19	31.79	36.68
	w/o DM	43.33	35.83	35.51	31.06	36.43	41.75	33.89	35.48	36.27	36.85
	Deepseek-V3+CIW	47.09	38.88	39.05	37.50	40.63	49.98	38.32	37.64	30.20	39.04
	w/o DM	46.40	35.15	37.19	37.65	39.15	48.89	41.03	36.69	31.99	39.65
	Qwen-plus+CIW	42.51	35.40	37.73	36.11	37.94	44.08	35.43	33.50	28.16	35.29
	w/o DM	40.77	35.73	35.66	35.23	36.85	42.16	32.91	31.85	28.45	33.84
	CIM-1.3B+CIW	42.13	31.76	31.47	30.37	33.93	42.94	33.10	36.59	36.32	37.24
	w/o DM	44.47	30.19	29.58	31.56	33.95	45.14	34.58	33.12	35.65	37.12
	CIM-6.7B+CIW	49.64	37.74	37.65	39.11	41.04	47.41	36.40	35.78	37.24	39.21
	w/o DM	44.92	35.91	36.34	34.90	38.02	48.16	37.76	35.02	37.51	39.61

TABLE X: Ablation study on ExeBench benchmark for Data Mapping(DM). (%)

Metric	Model	ExeBench 32-bit					ExeBench 64-bit				
		O0	O1	O2	O3	AVG	O0	O1	O2	O3	AVG
Re-com	GPT-4o+CIW	86.37	76.41	77.31	75.76	78.96	85.16	79.49	78.61	74.26	79.38
	w/o DM	71.08	61.27	61.57	59.09	63.25	80.75	72.47	72.50	75.60	75.33
	Deepseek-V3+CIW	83.69	77.46	76.39	78.79	79.08	91.05	83.43	81.67	80.43	84.15
	w/o DM	80.76	66.90	68.06	69.70	71.35	85.62	78.37	76.69	75.34	79.00
	Qwen-plus+CIW	70.32	67.96	71.76	73.23	70.82	76.10	71.07	73.61	67.29	72.02
	w/o DM	69.17	64.44	62.04	63.13	64.69	75.54	64.33	61.25	60.32	65.36
	CIM-1.3B+CIW	78.98	65.49	68.06	65.66	69.55	84.37	60.67	58.89	57.10	65.26
	w/o DM	72.36	65.85	64.35	62.12	66.17	80.18	60.39	58.33	58.45	64.34
	CIM-6.7B+CIW	85.99	72.54	68.06	69.19	73.94	86.86	62.64	60.56	59.25	67.33
	w/o DM	80.23	64.25	63.00	62.50	67.49	84.03	56.74	59.44	61.39	65.40
Re-exe	GPT-4o+CIW	26.62	14.44	10.65	6.57	14.57	36.35	9.83	7.50	6.17	14.96
	w/o DM	17.71	7.75	3.24	2.53	7.80	30.92	4.49	8.61	8.58	13.15
	Deepseek-V3+CIW	37.45	22.89	18.98	15.66	23.75	51.64	21.63	11.67	10.72	23.92
	w/o DM	25.86	11.97	7.41	4.55	12.45	40.88	13.20	12.20	8.58	18.71
	Qwen-plus+CIW	19.75	8.80	6.02	6.06	10.16	24.58	5.06	5.28	4.29	9.80
	w/o DM	9.55	4.93	2.31	3.03	4.96	21.86	3.37	4.88	4.83	8.73
	CIM-1.3B+CIW	49.04	25.35	18.06	16.16	27.15	55.83	23.31	14.72	13.40	26.82
	w/o DM	38.85	20.77	11.11	9.60	20.08	52.89	21.35	14.44	12.60	25.32
	CIM-6.7B+CIW	62.93	35.21	21.30	23.74	35.79	65.69	28.93	18.89	17.96	32.87
	w/o DM	52.03	23.58	14.98	14.42	26.25	60.93	22.47	17.50	17.43	29.58
ES	GPT-4o+CIW	47.00	32.73	31.92	33.23	36.22	49.48	33.18	33.21	32.23	37.03
	w/o DM	43.60	32.64	30.41	29.58	34.06	46.37	33.73	34.78	33.08	36.99
	Deepseek-V3+CIW	54.73	37.40	36.81	35.79	41.18	58.38	37.38	35.76	34.22	41.44
	w/o DM	53.29	36.78	37.69	35.23	40.75	51.13	34.79	34.53	32.94	38.35
	Qwen-plus+CIW	44.01	33.36	31.65	32.84	35.47	48.25	33.28	32.51	30.95	36.25
	w/o DM	41.33	30.66	30.42	30.48	33.22	45.18	33.21	32.87	31.41	35.67
	CIM-1.3B+CIW	67.99	50.00	42.51	39.43	49.98	68.43	47.89	39.43	38.28	48.51
	w/o DM	64.83	45.20	39.41	37.67	46.78	66.48	47.63	38.80	38.45	47.84
	CIM-6.7B+CIW	71.32	53.03	41.00	44.68	52.51	71.61	51.81	42.98	42.91	52.33
	w/o DM	68.66	51.85	42.12	40.52	50.79	70.38	48.24	44.18	42.28	51.27