Nova: Generative Language Models for Assem BLY CODE WITH HIERARCHICAL ATTENTION AND CON TRASTIVE LEARNING

Anonymous authors

Paper under double-blind review

Abstract

Binary code analysis is the foundation of crucial tasks in the security domain; thus building effective binary analysis techniques is more important than ever. Large language models (LLMs) although have brought impressive improvement to source code tasks, do not directly generalize to assembly code due to the unique challenges of assembly: (1) the low information density of assembly and (2) the diverse optimizations in assembly code. To overcome these challenges, this work proposes a *hierarchical attention* mechanism that builds attention summaries to capture the semantics more effectively, and designs *contrastive learning objectives* to train LLMs to learn assembly optimization. Equipped with these techniques, this work develops *Nova*, a generative LLM for assembly code. Nova outperforms existing techniques on binary code decompilation by up to 14.84 – 21.58% (absolute percentage point improvement) higher Pass@1 and Pass@10, and outperforms the latest binary code similarity detection techniques by up to 6.17% Recall@1, showing promising abilities on both assembly generation and understanding tasks.

1 INTRODUCTION

Binary code plays an irreplaceable role in the security domain, being the foundation of crucial tasks including vulnerability detection (Güler et al., 2019; Duan et al., 2020; Chen et al., 2022b), malware detection (Spensky et al., 2016; Aonzo et al., 2023; Xu et al., 2014), binary recovery (Su et al., 2024; Zhang et al., 2021; Chen et al., 2022c), and legacy software maintenance (Carbone et al., 2009; Carlini et al., 2015; Martin et al., 2010). For example, when performing tasks such as identifying attacks and malware, security analysts often only have access to assembly, i.e., the human-readable representation of binary code, which is extremely difficult to understand (Su et al., 2024; Zhang et al., 2021; Chen et al., 2022c). Thus, combined with the increasing sophistication of cybercrime that poses significant threats worldwide (e.g., cybercrime is predicted to cost the world \$10.5 trillion annually by 2025 (Sausalito, 2020)), effective binary analysis techniques are in high demand.



Figure 1: Example that shows the semantics and diverse optimizations of assembly code.

Large language models pre-trained on source code have brought improvement in various software development domains (Chen et al., 2022a; Liu et al., 2023; Chen et al., 2023; Le et al., 2022;
Jiang et al., 2023; Xia et al., 2023). However, these LLMs are not designed for or trained with assembly corpus, not achieving their full potential on binary code analysis tasks such as binary code similarity (Wang et al., 2022; Xu et al., 2023a), malware detection (Su et al., 2024), and binary code decompilation (Tan et al., 2024; Armengol-Estapé et al., 2024; Hosseini & Dolan-Gavitt, 2022).



Existing work applying LLMs on assembly code mainly piggybacks on encoder-style LLMs (Wang et al., 2022; Su et al., 2024; Xu et al., 2023a), unable to benefit from the more extensive pre-training, updated architectures, scaling of state-of-the-art generative LLMs. Other work using generative LLMs for decompilation shows a low unit test passing rate of the decompiled programs (Tan et al., 2024; Armengol-Estapé et al., 2024).

The challenges of leveraging generative LLMs for assembly code are twofold. First, compared to 060 source code, assembly code has a *lower information density*. A short source-code sequence maps 061 to an assembly-code sequence that is often several times longer. Thus, assembly semantics span 062 across a long sequence of tokens. Figure 1 (a) shows an example of a source code function that 063 compares two integers, while Figure 1 (b) shows its corresponding assembly code optimized with 064 00 flag. In the 00-optimized assembly code, the five instructions from 10: mov -0x8(%rbp),%rax to 1c: cmp %eax, %edx perform the checking whether the value of x is smaller than the value of y 065 (correspond to if (*(int*)x < *(int*)y) in the source code). A single assembly instruction alone 066 represents little meaningful semantics in the source code. It is the combinations of many instructions 067 and the *dependencies* between them represent the semantics. Such combinations of instructions are 068 long, which is hard for LLMs to learn. 069

Second, assembly code is diverse due to compiler optimization. The assembly code of the same source code function looks dramatically different with different compiler optimization. Figure 1
(c) shows the assembly of the same function compiled with 01 and 00 flags, which consists of a significantly different set of instructions. Such syntax diversity is hard for LLMs to learn, preventing LLMs from obtaining consistently good performances on differently optimized assembly code.

In this work, we develop Nova, a generative foundation LLM pre-trained for assembly code with two key novelties. First, to address the low-information-density and long-sequence challenge, we design a hierarchical self-attention, which contains three categories of attention at different levels of granularity: intra-instruction attention, preceding-instruction attention, and inter-instruction attention. The key insight is to build *attention summaries*, i.e., we create per-statement attention *labels*, which act as the summary of a statement. We then use preceding-instruction attention to capture semantics between a token and its preceding instruction label and use inter-instruction attention for long dependencies. Besides, we design *functionality contrastive learning* and *optimization contrastive learning* objectives to train Nova to learn the semantics behind the diverse syntax of assembly.

- 084 This work makes the following contributions:
 - We propose a novel hierarchical attention mechanism that captures the assembly's low-density semantics at three granularity levels.
 - We design contrastive learning objectives to train LLMs to learn assembly with diverse optimizations and encode assembly more efficiently.
 - We develop Nova, a generative foundation LLM with hierarchical attention and contrastive learning for assembly. Nova outperforms state-of-the-art on binary decompilation by up to 14.84 21.58% higher Pass@k and on binary similarity detection by up to 6.17% Recall@1.
 - 2 RELATED WORK

085

087

088

091

092 093

094

095 2.1 BINARY MODELS

Machine learning models are widely used in binary program analysis tasks. However, these models are typically designed for specific tasks such as binary code similarity detection (Pei et al., 2020; Xu et al., 2023a; Wang et al., 2022), variable name prediction (Chen et al., 2022c; Xu et al., 2023b; Zhang et al., 2021), binary code type inference (Pei et al., 2021), and are built from scratch.

Recent techniques have started to pre-train foundation LLMs for binaries. CodeArt (Su et al., 2024)
pre-trains encoder-style LLMs with a regularized attention design to better encode assembly code
semantics. SLaDe (Armengol-Estapé et al., 2024) trains BART (Lewis et al., 2019) models on
assembly. Meta LLMCompiler (Cummins et al., 2024) train CodeLlama models on LLVM IR to
optimize binary code. LLM4Decompile (Tan et al., 2024) trains DeepSeekCoder with assembly
for binary code decompilation. However, CodeArt does not generalize to generation tasks due to
its encoder architecture. LLMCompiler trained on LLVM IR cannot be effectively transferred to
assembly code. SLaDe and LLM4Decompile are limited in performance due to a lack of special

designs for assembly. In contrast, Nova proposes hierarchical attention and contrastive learning
 objectives, outperforming existing techniques on both understanding (binary code similarity detection)
 and generation (binary code decompilation) tasks.

112 2.2 LARGE SOURCE-CODE MODELS113

LLMs demonstrate promising results on many code-related tasks, such as code generation (Chen et al., 2022a; Liu et al., 2023; Chen et al., 2023; Le et al., 2022; Yue et al., 2021; Chen et al., 2021; Nijkamp et al., 2022; Fried et al., 2023; Rozière et al., 2023; Guo et al., 2024; Lozhkov et al., 2024; Hui et al., 2024), bug fixing (Jiang et al., 2023; Xia et al., 2023) and vulnerability fixing (Wu et al., 2023; Steenhoek et al., 2023; He & Vechev, 2023). The advances in using LLMs are attributed to the knowledge learned from massive source code and natural language text in their training datasets (Touvron et al., 2023; OpenAI, 2023). Nova is designed and trained for assembly, which has unique challenges such as low information density and diverse optimization.

121 122 123

2.3 ATTENTION MECHANISM

Standard self-attention is widely used in transformer architecture (Vaswani et al., 2017) to capture soft
dependencies between tokens in the input. Many special attention mechanisms have been designed
for better learning in various scenarios (Yang et al., 2016; Huang et al., 2024). LongCoder (Guo et al.,
2023) combines window attention and global attention (attention sink (Xiao et al., 2024)) to handle
long input of source code. We have shown that LongCoder's window attention is less effective than
Nova's on assembly code.

130 CAST (Shi et al., 2021) is a new neural architecture that splits the abstract syntax tree (AST) of source 131 code into subtrees, encodes the subtrees, and aggregates to the final encoding. PA-former (Chai & Li, 132 2024) is a new neural architecture that constructs source code as pyramid input based on their AST 133 structure and contains a pyramid attention mechanism to calculate the features in a hierarchical way. HierarchyNet (Nguyen et al., 2023) is a neural architecture that considers source code AST, data flow, 134 and control flow graphs. Similarly, it cannot be applied to assembly code. Different from CAST, 135 PA-former, and HierarchiyNet, Nova's attention design is for assembly code, is more lightweight and 136 can be plugged into any pre-trained generative LLM. 137

138 139

140

3 Approach

Figure 2 presents the overall approach of Nova. We build Nova on top of foundation models for source code (Guo et al., 2024) to utilize their source code and natural language generation ability. We first collect large assembly corpora (Section 3.1). Section 3.2 describes Nova's hierarchical attention design. With the collected assembly corpora, we then pretrain Nova with language modeling and contrastive learning objectives (Section 3.3). Then, we fine-tune Nova on two important downstream tasks, binary code decompilation, and binary code similarity detection (Sections 3.4 and 3.5), to prove Nova's effectiveness and benefits to the binary research domain.



3.1 DATA COLLECTION

Figure 2: Overview of developing Nova

Table 1: Statistics (number of C and X86-64 as-
sembly functions) of the pre-training datasets.

Origin	C Functions	O0	01	O2	O3	Total
AnghaBench	757.1K	743.1K	726.4K	718.7K	717.8K	3.7M
The-Stack	138.8K	125.1K	119.7K	116.9K	108.8K	609.3K

153 154

154 155

In this work, we focus on X86-64 assembly functions for C programs. Yet, Nova's approach is generalizable to other assembly languages such as ARM assembly.

We derive our X86-64 assembly functions dataset from two source code corpora: C functions in
 The-Stack (Li et al., 2023) and AnghaBench (da Silva et al., 2021). We compile the C programs into
 executables using gcc with different optimization levels (i.e., 00, 01, 02 and 03), strip the executables
 to remove debug information, and disassemble them into X86-64 assembly code using objdump. We

162 treat every function as a separate data point. Table 1 shows the number of C functions in the two 163 original datasets, and the number of X86-64 assembly functions we collected from them. 164

We perform certain normalization on the assembly functions: (1) removing all the "%" and comments, 165 (2) adding whitespace around ",", "(", ")", (3) converting all the hexadecimal numbers to decimal 166 numbers, and (4) replacing the address of each instruction with special labels (e.g., replacing "0" and "4" in Figure 1 (b) with "[INST-1]" and "[INST-2]") placing at the end of each instruction. More 168 details are in Appendix A.1. 169

HIERARCHICAL SELF-ATTENTION 3.2

167

170 171 172

173 174

175

176

177

178

204

Nova uses hierarchical self-attention that is specially designed to learn the *low-information-density* semantics in the *long* sequence of assembly code. Specifically, Nova learns the assembly code in an hierarchical way by providing a modified attention mask. Different from standard token-level attentions (Vaswani et al., 2017; Radford & Narasimhan, 2018; Radford et al., 2019; Brown et al., 2020), our hierarchical self-attention contains three categories at different levels of granularity.



Figure 3: Design of Nova's hierarchical attention for assembly code

205 (1) Intra-Instruction Attention: Due to the low information density in assembly, intra-instruction 206 attention is designed to capture the summary of every instruction, which is the standard causal attention but limited to tokens of each instruction (the yellow part in Figures 3) (a) and (b). Tokens 207 in different instructions have no attention weights. The "[INST]" label at the end of the instruction has 208 attention to all the tokens in the instruction and thus captures the semantics of the entire instruction 209 (e.g., "[INST-1]" captures the semantics of "mov eax, \$1"). 210

211 (2) Preceding-Instruction Attention: In addition to the local semantics of each instruction, the use 212 of assembly instructions (such as the choice of registers) depends on the context. For example, after 213 the first instruction "mov eax, \$1", the second instruction should not reuse "eax" to store another value "\$2" immediately. To capture such context, the preceding-instruction attention enables each token in 214 an instruction to have attention to the "[INST]" label of the preceding instruction (the light green 215 part in Figures 3 (a) and (b)).

(3) Inter-Instruction Attention: To understand function semantics (i.e., functionality), which lies in the dependencies across different instructions, the inter-instruction attention is designed to let the "[INST]" label of each instruction have attention to all the labels of previous instructions. For example, "[INST-4]" has attention to "[INST-1]", "[INST-2]", and "[INST-3]" (the dark green part in Figures 3 (a) and (b)). The inter-instruction attention is only enabled for "[INST]" labels, as they represent the semantics of each instruction.

222 To sum up, the hierarchical self-attention breaks the semantics of assembly code into three parts. The 223 intra-instruction attention captures the instruction summary, and the preceding-instruction attention 224 captures the context with the preceding instruction. The inter-instruction attention learns the long 225 dependencies across instructions on top of the "[INST]" labels that contain the instruction summary. 226 Figure 3 (c) also shows how Nova's hierarchical attention for assembly code is compatible with the 227 standard self-attention for text and source code. As existing LLMs have shown good performance on text and source code using the standard self-attention, we keep the standard causal attention 228 mask within and between any chunks of text or source code in the input (the light grey part shown in Figure 3 (c)). The attention from text or source code to assembly code (and vice versa) is only 230 allowed for the "[INST]" tokens as they are designed to contain the assembly instruction summaries. 231

233 3.3 CONTRASTIVE LEARNING

232

234

235

236

237

238

239

240

249

250 251

252

253

254

255

256

257

258

259

263 264 265 The syntax gap between assembly code and source code, and syntax diversity between differentlyoptimized assembly code make LLMs struggle to distinguish the semantics behind the syntax. Nova adopts contrastive learning technique (Gao et al., 2021) during pre-training to train LLMs to encode assembly code meaningfully w.r.t semantics. The standard pre-training objective is language modeling by minimizing the negative likelihood of code in the pre-training corpus (Radford & Narasimhan, 2018), notated as L_{lm} . In addition, Nova is pre-trained with two new objectives, L_{fcl} for functionality contrastive learning and L_{ocl} for optimization contrastive learning.



Figure 4: Design of functionality and optimization contrastive learning (CL). "asm" denotes assembly.

Functionality CL: Functionality CL trains Nova to focus more on the functionalities of assembly code rather than the syntax. Code with the same functionality (assemblies from the same source code), should be encoded closer in the latent space. For instance, in Figure 4 (a), embeddings of source and assembly code of function "cmp" are closer to each other, and the same for function "sort". Nova is designed and implemented on decoder-only generative LLM, and we refer the hidden states from the last transformer layer as embedding. For source code, we use the average of each token's embedding as the source code function's embedding. For assembly, we use the average of all the "[INST]" tokens' embedding as the embedding of the assembly function, as each "[INST]" token is supposed to capture the semantics of that instruction by the design of hierarchical self-attention.

Let e_f^s be the embedding of function f in s form (s = -1 for source code, and $s \in [0, 1, 2, 3]$ for 00 to 03 optimized assembly). For simplicity, let S = [-1, 0, 1, 2, 3] be the domain of s. Functionality CL optimizes Nova with the constraint:

$$\forall f_i \in F, \ \max_{s,t \in S} (d(e_{f_i}^s, e_{f_i}^t)) < \min_{\substack{s,t \in S \\ f_j \neq f_i \in F}} (d(e_{f_i}^s, e_{f_j}^t))$$

, where d calculates the l_2 distance between two embeddings and F is the full set of functions in the training corpus.

Such constraints can be trained by optimizing the embeddings of a batch of functions, each function in two different forms. For the example in Figure 4 (b), there are two forms (source code and 00

assembly) of four functions. Once Nova encodes the batch of source code and assembly functions, we calculate the distance matrix $\{D_{ij}\}_{f_i,f_j \in F} = \{d(e_{f_i}^s, e_{f_j}^t)\}$, and minimize the loss:

273 274

275 276

294

295 296

297

305

306

$$L_{fcl} = -\log \sum_{s,t \in S} \sum_{f_i \in F} \left(1 - \frac{\exp\left(d(e_{f_i}^s, e_{f_i}^t)\right)}{\sum_{f_j \in F} \exp\left(d(e_{f_i}^s, e_{f_j}^t)\right)} \right)$$

This objective minimizes the distance between embeddings for the same function, i.e., the diagonal in the distance matrix.

Optimization CL: LLMs can be confused if being asked to directly connect a source code function to its 03-optimized assembly, due to their dramatically different syntax. Such a huge gap can be filled by learning how the source code is transformed to 00, 01, 02 and eventually to 03 assembly, as the optimization levels are *ordered*.

Higher-level optimization applies a super-set of optimization rules compared to lower-level optimization. Nova learns such order with the optimization CL objective, encoding differently-optimized assembly code orderly. Optimization CL optimizes Nova with the constraint: the more optimizations applied, the larger the difference between embeddings of optimized and unoptimized code. For instance, Figure 4 (c) and (d) illustrate that for the same function "cmp", the distance between source code and assembly increases when the optimization level increases. Formally, optimization CL minimizes the following loss:

$$L_{ocl} = \sum_{f \in F} \sum_{s < t_1 < t_2 \in S} max \left(0, d(e_f^s, e_f^{t_1}) - d(e_f^s, e_f^{t_2}) \right)$$

Overall, the final training loss combines the three: $L = L_{lm} + \lambda (L_{fcl} + L_{ocl})$, where λ is set to 0.1 to balance the losses in this work.

3.4 TASK 1: BINARY CODE DECOMPILATION

Binary code decompilation (BCD) helps developers to understand binary code by recovering binary code into more readable high-level source code (e.g., C programs) (Fu et al., 2019; Liang et al., 2021; Armengol-Estapé et al., 2024; Tan et al., 2024). The input to the model for BCD is formatted as an instruction prompt (notated by p): # This is the assembly code with {opt} optimization: {asm}, where "opt" is the optimization-level applied to the assembly and "asm" is the assembly code to decompile. Nova is fine-tuned to generate the expected source code function src following the instruction prompt. The fine-tuning objective is minimizing the loss: $L_{bcd} = -\log P(src|p)$.

3.5 TASK 2: BINARY CODE SIMILARITY DETECTION

Binary code similarity detection (BCSD) aims to measure the similarity between two binary code
snippets (Wang et al., 2022; Su et al., 2024), which is the foundation of various applications such as
plagiarism detection (Luo et al., 2014; Sæbjørnsen et al., 2009) and vulnerability detection (David &
Yahav, 2014; David et al., 2018; 2017; 2016).

A widely used setting is taking a query assembly of the function f^q that is compiled with one optimization level (denoted by *s*), and a pool of candidate assembly of *K* (e.g., 50, 100, etc.) functions (notated by f_i^p , $1 \le i \le K$) compiled with a different optimization level (denoted by $t \ne s$). There exists a unique candidate assembly coming from the same source code as the query ($\exists ! 1 \le i \le K, f_i^p = f^q$, called the positive candidate). Nova is fine-tuned to encode these binaries, so that the positive candidate has the highest similarity with the query assembly among the pool. The learning objective is as follows:

318 319 320

321 322

$$L_{BCSD} = -\log \sum_{1 \le i \le K}^{f^q := f_i^p} \left(1 - \frac{\exp\left(d(e_{f^q}^s, e_{f_i^p}^t)\right)}{\sum_{1 \le j \le K} \exp\left(d(e_{f^q}^s, e_{f_j^p}^t)\right)} \right)$$

We follow previous work (Su et al., 2024) to let s be 00-assembly and t be 03-assembly, which is the hardest setting.

324 4 EXPERIMENTAL SETUP

This section describes the setup of pre-training and fine-tuning of Nova, as well as the existing baselines we compare Nova with, and the evaluation metrics we used in the two downstream tasks. Appendix A.2 contains additional details such as training hyper-parameters.

4.1 PRE-TRAINING

We use the C and X86-64 assembly functions collected from AnghaBench and The-Stack for pretraining. We pre-train Nova starting from DeepSeek-Coder (Guo et al., 2024), and the hierarchical attention is applied on half of the attention heads to balance between its effectiveness and the existing knowledge in the standard attention layers (Justified in Appendix A.3). Nova is pre-trained with language modeling for one epoch, followed by contrastive learning objectives for another epoch.

337 338

326

327

328

329 330

331

4.2 FINE-TUNING FOR BINARY CODE DECOMPILATION

Training Data: We sample (due to computation resource limitation) 2.16M assembly-to-source-code pairs (0.338B tokens) from the pre-training corpus to build the BCD fine-tuning data.

Test Data: We use HumanEval-Decompile (Tan et al., 2024) as the test benchmark, which was not used in training. HumanEval-Decompile is derived from the C language adaptation of the HumanEval (Chen et al., 2021) benchmark and contains 164 C functions, each compiled with 00 – 03 optimization flags and disassembled into X86-64 assembly.

Baselines: Nova is compared with existing SOTA LLM4Decompile (Tan et al., 2024), Meta LLMCompiler (Cummins et al., 2024), and other general code LLMs (CodeLlama (Rozière et al., 2023),
StarCoder2 (Lozhkov et al., 2024), DeepSeekCoder (Guo et al., 2024), Qwen2.5-Coder (Hui et al.,
2024), Granite Mishra et al. (2024), LongCoder (Guo et al., 2023), GPT-3.5-Turbo, and GPT-40).
LLM4Decompile trains DeepSeekCoder using the same AnghaBench corpus for binary decompilation. Meta LLMCompiler trains CodeLlama models using LLVM IRs, X86, and ARM assembly
code to optimize and translate binary code.

Evaluation: We let each model sample 20 decompilations per assembly function, using the temperature of 0.2 and top_p of 0.95 (Chen et al., 2021). Except for LLM4Decompile and Nova that are fine-tuned for binary code decompilation, we provide all other baselines with three-shot examples for few-shot learning (Brown et al., 2020).

357 358

359

4.3 FINE-TUNING FOR BINARY CODE SIMILARITY DETECTION

Training Data: To compare Nova with existing works on BCSD fairly (Wang et al., 2022; Su et al., 2024), we use the same dataset, BinaryCorp-3M (Wang et al., 2022), as the fine-tuning data for BCSD, which contains the 00 and 03 assembly of 224,606 functions.

Test Data: Following existing work (Su et al., 2024; Xu et al., 2023a), we use real-world benchmarks,
 Binutils, Curl, ImageMagick, SQLite, OpenSSL, and Putty, as the test benchmarks, which are
 nonexistent in the training data.

Baselines: Nova is compared with jTrans (Wang et al., 2022), DiEmph (Xu et al., 2023a) and
CodeArt (Su et al., 2024). jTrans is a Transformer (Vaswani et al., 2017) encoder trained on
binaries with masked token prediction and jump target prediction tasks. DiEmph uses an instruction
deemphasis technique to prevent the model from learning instruction distribution biases introduced by
compilers. CodeArt proposes a regularized attention mask for encoder models to capture instructional
semantics and data dependencies.

Evaluation: We randomly sample K (K = 50, 100, 200, 500) source code functions from each project, compile them into binaries with 00 and 03 optimization flags, and disassemble them into X86-64 assemblies. BCSD techniques encode these assemblies into embeddings (Nova uses the average last-layer hidden states of all the "[INST]" tokens in an assembly as its embedding). Then each 00 assembly is used as the query to calculate their similarity with the K 03 candidate assemblies (using cosine similarity). Metric Recall@1 is the ratio of queries for which the candidate from the same source code has the highest similarity among all the candidates.

5 Results

5.1 BINARY CODE DECOMPILATION

Comparison with SOTA Techniques: Table 2 shows the Pass@1 and Pass@10 of the decompiled code from assemblies on HumanEval-Decompile. The results are grouped by optimization level (i.e., the benchmark contains 164 assemblies of each optimization level), and the average is also reported.

Table 2: Nova's Pass@K and comparison with existing techniques on HumanEval-Decompile.

T 1 :			Pass@1					Pass@10		
Techniques	00	01	02	O3	Avg.	00	01	O2	O3	Avg.
CodeLlama-7B	6.95	3.81	4.54	3.78	4.77	8.53	5.97	7.34	5.17	6.75
StarCoder2-7B	6.31	4.33	5.64	5.95	5.56	8.77	5.18	6.09	7.17	6.80
DeepSeekCoder-7B	9.63	7.56	7.41	6.68	7.82	13.60	11.38	11.52	9.44	11.49
Qwen-2.5-Coder-7B	4.76	5.79	5.58	5.27	5.35	6.34	7.69	6.56	5.79	6.60
LLMCompiler-7B	5.95	5.85	5.55	5.82	5.79	7.01	7.31	7.47	7.01	7.20
GPT-3.5-Turbo	7.41	6.13	4.33	3.90	5.44	9.56	8.38	6.23	5.12	7.32
GPT-40	21.34	18.29	14.48	13.05	16.79	29.94	26.74	21.42	19.88	24.50
LLM4Decompile-1B	15.30	8.26	9.36	8.38	10.33	21.79	15.23	16.17	13.70	16.72
Nova-1B	37.53	21.71	22.68	18.75	25.17	49.38	34.84	36.95	32.03	38.30
LLM4Decompile-6B	29.97	19.05	20.46	18.32	21.95	40.40	27.75	28.85	28.51	31.38
Nova-6B	48.78	30.58	30.85	27.23	34.36	57.47	47.45	43.03	39.68	46.91

Overall, Nova's Pass@1 and Pass@10 are higher than all SOTA binary decompilation techniques
 and general LLMs with even smaller model sizes. Specifically, for each optimization level, Nova
 consistently decompiles more assemblies into source code correctly than the rest of the compared
 techniques. Note that Meta LLMCompiler is mainly designed for LLVM IR optimization, and thus is
 still incapable of assembly code decompilation.

With the same model size, Nova-1B outperforms LLM4Decompile-1B with a 14.84% higher averaged
Pass@1, and a 21.58% higher Pass@10. Nova-6B outperforms LLM4Decompile-6B with a 12.41%
higher averaged Pass@1, and a 15.53% higher Pass@10. When compared with GPT-4o, an order of
magnitude larger model, Nova-1B produces an 8.38% higher Pass@1 and 13.80% higher Pass@10.
Examples of Nova's correct decompilation are provided in Appendix A.5.

Comparison with Techniques Handling Long Input: Nova's hierarchical attention design targets to address the low information density and long input challenge of assembly code. There are other techniques that handles long input challenges in text and source code, with Granite-3B-Code-Base-128K and LongCoder being the most related ones. Granite trains LLM on repository-level long inputs, which is an orthogonal approach with Nova's approach (hierarchical attention and contrastive learning). We train Granite-3B-Code-128K with Nova's approach, and Table 3 shows that Nova's approach brings improvement to Granite over standard fine-tuning even if it has already been trained with long code data.

Tashairman			Pass@1					Pass@10		
Techniques	00	01	O2	O3	Avg.	00	01	02	O3	Avg.
Granite (3B-Code-128K)	5.91	3.78	5.09	5.52	5.08	8.19	5.16	6.56	7.15	6.76
Granite + Standard Fine-Tuning	20.88	13.54	11.37	10.09	13.97	30.05	19.77	18.31	15.77	20.98
Granite + Nova's Approaches	31.04	14.57	14.70	13.66	18.49	39.57	21.23	21.77	19.82	25.60

Table 3: Nova's approach brings improvement to LLM that has been trained with long input data.

Teshniswaa			Pass@1					Pass@10		
Techniques	00	O1	02	O3	Avg.	00	01	02	O3	Avg.
Nova-1B (using LongCoder's Attention) Nova-1B	34.59 37.53	19.07 21.71	19.72 22.68	17.34 18.75	22.68 25.17	42.19 49.38	32.37 34.84	32.86 36.95	29.04 32.03	34.12 38.30

431 LongCoder combines window attention and global attention to learn long code input. We compare LongCoder's attention design with Nova's hierarchical attention design by replacing the hierarchical

Table 5: Ablation stu	ıdy	of Nova-1B on	HumanEval-Decom	pile.
				r · ·

Techniques			Pass@1					Pass@10		
reeninques	00	01	02	03	Avg.	00	01	02	O3	Avg.
LLM4Decompile-1B	15.30	8.26	9.36	8.38	10.33	21.79	15.23	16.17	13.70	16.72
Nova_ <i>CL</i> - <i>HA</i> Nova_ <i>HA</i>	20.73 30.58	16.16 19.88	15.03 20.58	11.19 16.40	15.78 21.86	33.55 44.75	28.12 33.13	26.96 33.31	21.01 29.82	27.41 35.25

attention of Nova-1B with LongCoder's attention design. Table 4 shows that Nova's hierarchical attention is more effective in learning assembly code. Nova's attention design considers the instruction-level local semantics and dependencies between different instructions, which fits better than fix-sized window attention to assembly code.

Ablation Study: We conduct an ablation study by comparing Nova-1B with the following models to show the effectiveness of contrastive learning objectives and hierarchical attention design:

• Nova_{-CL-HA}: Removing contrastive learning and hierarchical self-attention. This is simply training DeepSeekCoder-1.3B on the assembly corpus using language modeling. This can be viewed as our reproduction (retrain) of LLM4Decompile-1B using the same amount of data.

Nova_{-HA}: Removing the hierarchical self-attention, training DeepSeekCoder-1.3B on the assembly corpus using both the language modeling and contrastive learning objectives.

Table 5 shows the results of the ablation study. Nova $_{CL-HA}$ produces an average Pass@1 of 15.78% and Pass@10 of 27.41%. With additional contrastive learning objectives, Nova $_{HA}$ improves the Pass@1 on all optimization levels over Nova $_{CL-HA}$, showing a higher averaged Pass@1 and Pass@10. Further applying the hierarchical self-attention on Nova $_{HA}$ boosts the overall Pass@1 from 21.86% to 25.17%, and Pass@10 from 35.25% to 38.30%.

Table 6: Recall@1 on BCSD with K = 50

Table 7: Recall@1 on BCSD with K = 100

Benchmarks	jTrans	DiEmph	CodeArt	Nova-1B	Nova-6B
Binutils	0.68	0.80	0.84	0.87	0.89
Curl	0.72	0.84	0.86	0.89	0.94
ImageMagick	0.53	0.71	0.78	0.86	0.90
SQLite	0.73	0.79	0.78	0.77	0.78
OpenSSL	0.70	0.83	0.88	0.90	0.92
Putty	0.63	0.72	0.69	$\overline{0.72}$	0.71
Avg.	0.67	0.78	0.81	0.84	0.86

Benchmarks	fTrans	DiEmph	CodeArt	Nova-1B	Nova-6B
Binutils	0.60	0.63	0.81	0.79	0.79
Curl	0.63	0.80	0.82	0.86	0.88
ImageMagick	0.54	0.71	0.76	0.79	0.81
SQLite	0.62	0.72	0.74	0.73	0.72
OpenSSL	0.60	0.80	0.87	0.88	0.90
Putty	0.58	0.64	0.64	<u>0.65</u>	0.64
Avg.	0.60	0.72	0.77	0.78	0.79

Table 8: Recall@1 on BCSD with K = 200

Benchmarks	jTrans	DiEmph	CodeArt	Nova-1B	Nova-6B
Binutils	0.51	0.64	0.74	0.73	0.73
Curl	0.57	0.77	0.78	0.83	0.84
ImageMagick	0.39	0.51	0.67	0.73	0.75
SQLite	0.56	0.65	0.68	0.68	0.69
OpenSSL	0.54	0.71	$\widetilde{0.82}$	$\widetilde{0.84}$	0.88
Putty	0.49	<u>0.58</u>	0.55	0.55	0.58
Avg.	0.51	0.64	0.71	0.73	0.75

Benchmarks	jTrans	DiEmph	CodeArt	Nova-1B	Nova-6B
Binutils	0.40	0.57	0.70	0.65	0.67
Curl	0.43	0.62	0.69	0.73	0.76
ImageMagick	0.25	0.42	0.58	0.61	0.65
SQLite	0.43	0.59	0.62	0.59	0.62
OpenSSL	0.43	0.61	0.76	0.78	0.82
Putty	0.38	0.50	0.49	0.47	0.51
Avg.	0.39	0.55	0.64	0.64	0.67

5.2 BINARY CODE SIMILARITY DETECTION

Tables 6, 7, 8 and 9 show the Recall@1 of Nova and existing BCSD techniques with pool size K of 50, 100, 200 and 500 on the six benchmarks. <u>Underline</u> indicates the best in each benchmark, while wave denotes the tied best (we only mark Nova-1B for clearer illustration).

Overall, Tables 6, 7, 8 and 9 show that *on average, Nova-1B and Nova-6B achieve the highest Recall@1 (in bold) under all four settings of K*. Nova-6B further outperforms Nova-1B and achieves
the highest averaged Recall@1 under all four settings, ranking the ground-truth of 5%, 2%, 4%, and
3% more queries the most similar correspondingly compared to CodeArt. Nova-1B consistently
outperforms existing techniques with higher Recall@1 when K is 50, 100, and 200, meaning it

correctly ranks ground-truth of 3%, 1%, and 2% more queries as the most similar. Under the setting of K = 500, Nova-1B ties with CodeArt with the same highest Recall@1. When looking into each individual benchmark, Nova-1B always wins on the most benchmarks under different settings of pool size K. For instance, Nova-1B wins on four benchmarks while DiEmph only wins on SQLite when K = 50. We also conduct an ablation study on BCSD in Appendix A.6.

5.3 ANALYTIC EXPERIMENTS: HOW ARE NOVA'S EMBEDDINGS BETTER?

We use the widely-used t-SNE (van der Maaten & Hinton, 2008) to analyze and visualize highdimensional embeddings. We randomly sample seven coding problems from HumenEval-Decompile (task_id 19, 32, 34, 63, 119, 128, 143), encode the 00 - 03 assemblies by Nova_{-CL-HA} and Nova-1B. Figure 5 shows the embeddings that are visualized under the first two principal components. Each color represents one task, and 00 - 03 assemblies are marked by \bigcirc , \bigtriangledown , \triangle , and \Box .

Compared with Nova_{-CL-HA} (Figure 5 (a)), Nova_{<math>-HA} (Figure 5 (b)) including contrastive learning objectives in the pre-training, can separate the embeddings of assemblies with different functionalities better. Nova_{-HA} clearly encode "Task 143" (orange points) away from the others. Yet, Nova's(Figure 5 (c)) embeddings group the assemblies by functionalities more precisely than Nova_{<math>-HA},suggesting that hierarchical attention enhances the training of contrastive learning objectives to learnmore effective encoding. Visualization using a different approach, PCA, is shown in Appendix A.7.Analytic experiments on Nova's hierarchical attention is shown in Appendix A.8.</sub></sub></sub></sub>



Figure 5: t-SNE analysis of embeddings calculated by Nova $_{-CL-HA}$, Nova $_{-HA}$, and Nova.

6 LIMITATIONS

One limitation is that Nova is X86-specific, as we only collect X86 assembly corpus for pre-training. 521 This design choice is mainly affected by two facts: (1) X86 assembly is used and explored in a 522 wide range of binary tasks (Wang et al., 2022; Su et al., 2024; Xu et al., 2023a; Chen et al., 2022c) 523 compared to other assembly languages, and (2) computation limitations. However, the proposed 524 techniques are independent of X86 assembly. Low information density and compiler optimization are 525 the common challenges of most assembly languages such as X86, ARM, and MIPS. The proposed 526 techniques can be applied to the future development of multi-lingual assembly LLMs. Another 527 potential limitation is the scale of models. We develop Nova-1B and Nova-6B. These two LLMs 528 show impressive ability in assembly code decompilation and encoding. There might be a potential benefit of developing larger Nova models. However, due to the computing resources limitation, we 529 are unable to explore that in this work. 530

531 532

533

491 492

493

506

507

509

510

511

512

513

514

515

516 517 518

519 520

7 CONCLUSION

This work develops Nova, a generative foundation LLM for assembly code, which incorporates
two key novelties (hierarchical attention and contrastive learning objectives) to address the unique
challenges of assembly code. Evaluation on downstream tasks shows the effectiveness of Nova, which
outperforms existing techniques on binary code decompilation by up to 146.54% and outperforms
the latest binary code similarity detection techniques by up to 6.17%. We expect our hierarchical
attention and contrastive learning techniques to benefit source code and natural language foundation
models, which remains as future work.

540 REFERENCES

548

576

- Simone Aonzo, Yufei Han, Alessandro Mantovani, and Davide Balzarotti. Humans vs. Machines
 in Malware Classification. In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 1145–1162, Anaheim, CA, 2023. USENIX Association. ISBN 978-1-939133-37-3. URL https://www.usenix.org/conference/usenixsecurity23/presentation/aonzo.
- Jordi Armengol-Estapé, Jackson Woodruff, Chris Cummins, and Michael F. P. O'Boyle. SLaDe: A
 Portable Small Language Model Decompiler for Optimized Assembly, 2024.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. *CoRR*, abs/2005.14165, 2020. URL https://arxiv.org/abs/2005.14165.
- Martim Carbone, Weidong Cui, Long Lu, Wenke Lee, Marcus Peinado, and Xuxian Jiang. Mapping kernel objects to enable systematic integrity checking. In *Proceedings of the 16th ACM Conference on Computer and Communications Security*, CCS '09, pp. 555–565, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605588940. doi: 10.1145/1653662.1653729. URL https://doi.org/10.1145/1653662.1653729.
- 561 Nicholas Carlini, Antonio Barresi, Mathias Payer, David Wagner, and Thomas R. Gross. Control 562 Flow Bending: On the Effectiveness of Control-Flow Integrity. In 24th USENIX Security Sym 563 posium (USENIX Security 15), pp. 161–176, Washington, D.C., 2015. USENIX Association.
 564 ISBN 978-1-939133-11-3. URL https://www.usenix.org/conference/usenixsecurity15/
 565 technical-sessions/presentation/carlini.
- Lei Chai and Ming Li. Pyramid attention for source code summarization. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.
- Angelica Chen, Jérémy Scheurer, Tomasz Korbak, Jon Ander Campos, Jun Shern Chan, Samuel R.
 Bowman, Kyunghyun Cho, and Ethan Perez. Improving Code Generation by Training with Natural
 Language Feedback, 2023.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen.
 CodeT: Code Generation with Generated Tests, 2022a. URL https://arxiv.org/abs/2207.
 10397.
- Ligeng Chen, Zhongling He, Hao Wu, Fengyuan Xu, Yi Qian, and Bing Mao. DIComP: Lightweight
 Data-Driven Inference of Binary Compiler Provenance with High Accuracy. In 2022 IEEE *International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 112–122, 2022b. doi: 10.1109/SANER53432.2022.00025.
- 581 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared 582 Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, 583 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, 584 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios 585 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, 586 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, 588 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob 589 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating 590 Large Language Models Trained on Code. CoRR, abs/2107.03374, 2021. URL https://arxiv. org/abs/2107.03374. 592
- Qibin Chen, Jeremy Lacomis, Edward J. Schwartz, Claire Le Goues, Graham Neubig, and Bogdan Vasilescu. Augmenting Decompiler Output with Learned Variable Names and Types. In *31st*

607

618

627

631

632

633 634

635

636 637

 USENIX Security Symposium (USENIX Security 22), pp. 4327–4343, Boston, MA, aug 2022c.
 USENIX Association. ISBN 978-1-939133-31-1. URL https://www.usenix.org/conference/ usenixsecurity22/presentation/chen-qibin.

- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What Does BERT Look at? An Analysis of BERT's Attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes (eds.), *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 276–286, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828. URL https://aclanthology.org/W19-4828.
- Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Roziere, Jonas Gehring, Gabriel Synnaeve, and Hugh Leather. Meta Large Language Model Compiler: Foundation Models of Compiler Optimization, 2024. URL https://arxiv.org/abs/2407.02524.
- Anderson Faustino da Silva, Bruno Conde Kind, José Wesley de Souza Magalhães, Jerônimo Nunes Rocha, Breno Campos Ferreira Guimarães, and Fernando Magno Quinão Pereira. ANGHABENCH: A Suite with One Million Compilable C Benchmarks for Code-Size Reduction. In 2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO), pp. 378–390, 2021. doi: 10.1109/CGO51591.2021.9370322.
- Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David and Eran Yahav. Tracelet-Based Code Search in Executables. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI
 Yaniv David And Side Search in *PLOI (145/2594291)*. 2594343. URL https://doi.org/10.1145/2594291. 2594343.
- Yaniv David, Nimrod Partush, and Eran Yahav. Statistical Similarity of Binaries. In *Proceedings* of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '16, pp. 266–280, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342612. doi: 10.1145/2908080.2908126. URL https://doi.org/10.1145/2908080. 2908126.
- Yaniv David, Nimrod Partush, and Eran Yahav. Similarity of Binaries through Re-Optimization.
 SIGPLAN Not., 52(6):79–94, jun 2017. ISSN 0362-1340. doi: 10.1145/3140587.3062387. URL https://doi.org/10.1145/3140587.3062387.
- Yaniv David, Nimrod Partush, and Eran Yahav. FirmUp: Precise Static Detection of Common Vulnerabilities in Firmware. *SIGPLAN Not.*, 53(2):392–404, mar 2018. ISSN 0362-1340. doi: 10.1145/3296957.3177157. URL https://doi.org/10.1145/3296957.3177157.
 - Yue Duan, Xuezixiang Li, Jinghan Wang, and Heng Yin. DeepBinDiff: Learning Program-Wide Code Representations for Binary Diffing. 01 2020. doi: 10.14722/ndss.2020.24311.
 - Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen tau Yih, Luke Zettlemoyer, and Mike Lewis. InCoder: A Generative Model for Code Infilling and Synthesis, 2023.
- Cheng Fu, Huili Chen, Haolan Liu, Xinyun Chen, Yuandong Tian, Farinaz Koushanfar, and Jishen Zhao. Coda: An End-to-End Neural Program Decompiler. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/ 093b60fd0557804c8ba0cbf1453da22f-Paper.pdf.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple Contrastive Learning of Sentence
 Embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6894–6910, Online and Punta Cana, Dominican Republic, November 2021.
 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL https:
 //aclanthology.org/2021.emnlp-main.552.

648 649	Emre Güler, Cornelius Aschermann, Ali Abbasi, and Thorsten Holz. AntiFuzz: Impeding Fuzzing Audits of Binary Executables. In 28th USENIX Security Symposium (USENIX Security 19), pp.
650 651	1931–1947, Santa Clara, CA, 2019. USENIX Association. ISBN 978-1-939133-06-9. URL https://www.usenix.org/conference/usenixsecurity19/presentation/guler.
652	
653	Daya Guo, Canwen Xu, Nan Duan, Jian Yin, and Julian McAuley. LongCoder: a long-range pre-
654	trained language model for code completion. In <i>Proceedings of the 40th International Conference</i>
655	on Machine Learning, ICMIL 25. JMILK.org, 2025.
656	Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
657	Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. DeepSeek-Coder: When the
658	Large Language Model Meets Programming – The Rise of Code Intelligence, 2024.
659	lingvuan He and Martin Vechey Large Language Models for Code: Security Hardening and
660 661	Adversarial Testing. CoRR, abs/2302.05319, 2023. URL https://arxiv.org/abs/2302.05319.
662	Iman Hosseini and Brendan Dolan-Gavitt Beyond the C: Retargetable Decompilation using Neural
662	Machine Translation. In Proceedings 2022 Workshop on Binary Analysis Research, BAR 2022.
664	Internet Society, 2022. doi: 10.14722/bar.2022.23009. URL http://dx.doi.org/10.14722/
665	bar.2022.23009.
666	Vennen Herry Lizzari, V., Lerry Lei, Zim Ling, Tealer Chen, Zang Li, Ven Ven Visering Me
667	Lijuan Vang, Hao Chen, Shupeng Li, and Penghao Zhao. Advancing Transformer Architecture in
668	Long-Context Large Language Models: A Comprehensive Survey 2024 URL https://arxiv
669	org/abs/2311.12351.
670	
671	Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,
672	Bowen Yu, Kai Dang, An Yang, Kui Men, Fel Huang, AingZhang Ken, Auancheng Ken, Jingren Zhou, and Junyang Lin. Owan? 5 Coder Technical Peport. arXiv praprint arXiv:2400.12186
673	2024
674	
675	Nan Jiang, Kevin Liu, Thibaud Lutellier, and Lin Tan. Impact of Code Language Models on
676	Automated Program Repair. In Proceedings of the 45th International Conference on Software
677	Engineering, ICSE 25, pp. 1450–1442. IEEE Press, 2025. ISBN 9781005457019. doi: 10.1109/ ICSE48610.2023.00125_UBL https://doi.org/10.1109/ICSE48610.2023.00125
678	ICSE48019.2023.00125. ORE https://doi.org/10.1103/1CSE48019.2025.00125.
679	Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven C.H. Hoi. CodeRL:
680	Mastering Code Generation through Pretrained Models and Deep Reinforcement Learning. arXiv
681	<i>preprint</i> , abs/2207.01780, 2022.
682	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
683	Ves Stoyanov, and Luke Zettlemoyer. BART: Denoising Sequence-to-Sequence Pre-training for
684	Natural Language Generation, Translation, and Comprehension, 2019.
COO	Raymond Li Loubna Ben Allal Yangtian Zi Niklas Muennighoff Denis Kocetkov Chenghao Mou
607	Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Oian Liu, Evgenii Zheltonozhskii, Terry Yue
699	Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro,
620	Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar
690	Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason
691	Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang,
692	Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas,
693	waxini Kunakov, redor Zhdanov, Manuel Komero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger Hailey Schoelkonf Jan Ebert Tri Dao Mayank Michra Alay Cu Jennifer Pohinson
694	Carolyn Jane Anderson Brendan Dolan-Gavitt Danish Contractor Siva Reddy Daniel Fried
695	Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis. Sean Hughes. Thomas Wolf. Ariun
696	Guha, Leandro von Werra, and Harm de Vries. StarCoder: may the source be with you!, 2023.
697	Duisons Lions Ving Coo Daiwai Un and Kai Chan Nautanna attention based as at 1
698	Cybersecurity A(1):5 2021 ISSN 2523 3246 doi: 10.1186/s42400.021.00070.0. LIDL https://
699	//doi org/10 1186/s42400-021-00070-0
700	,, dollo. ₀ , lot 100, oll loo oll oooro ol
701	liate Liu Yigin Zhu Kaiwen Xiao Ojang Fu Xiao Han Wei Yang and Deheng Ve RITE

701 Jiate Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: Reinforcement Learning from Unit Test Feedback, 2023.

- 702 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane 703 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, 704 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, 705 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, 706 Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri 708 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa 710 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, 711 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. StarCoder 2 and The Stack 712 v2: The Next Generation, 2024. URL https://arxiv.org/abs/2402.19173. 713
- Lannan Luo, Jiang Ming, Dinghao Wu, Peng Liu, and Sencun Zhu. Semantics-Based ObfuscationResilient Binary Code Similarity Comparison with Applications to Software Plagiarism Detection.
 In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, FSE 2014, pp. 389–400, New York, NY, USA, 2014. Association for Computing
 Machinery. ISBN 9781450330565. doi: 10.1145/2635868.2635900. URL https://doi.org/10.
 1145/2635868.2635900.
- Jean-Phillipe Martin, Michael Hicks, Manuel Costa, Periklis Akritidis, and Miguel Castro. Dynamically checking ownership policies in concurrent c/c++ programs. *SIGPLAN Not.*, 45 (1):457–470, jan 2010. ISSN 0362-1340. doi: 10.1145/1707801.1706351. URL https://doi.org/10.1145/1707801.1706351.
- 724 Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, 725 Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, Manish Sethi, Xuan-726 Hong Dang, Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew Coleman, Matthew White, Mark 727 Lewis, Raju Pavuluri, Yan Koyfman, Boris Lublinsky, Maximilien de Bayser, Ibrahim Abdelaziz, 728 Kinjal Basu, Mayank Agarwal, Yi Zhou, Chris Johnson, Aanchal Goyal, Hima Patel, Yousaf 729 Shah, Petros Zerfos, Heiko Ludwig, Asim Munawar, Maxwell Crouse, Pavan Kapanipathi, Shweta 730 Salaria, Bob Calio, Sophia Wen, Seetharami Seelam, Brian Belgodere, Carlos Fonseca, Amith 731 Singhee, Nirmit Desai, David D. Cox, Ruchir Puri, and Rameswar Panda. Granite Code Models: A Family of Open Foundation Models for Code Intelligence, 2024. URL https://arxiv.org/ 732 abs/2405.04324. 733
- Minh Huynh Nguyen, Nghi D. Q. Bui, Truong Son Hy, Long Tran-Thanh, and Tien N. Nguyen. HierarchyNet: Learning to Summarize Source Code with Heterogeneous Representations, 2023. URL https://arxiv.org/abs/2205.15479.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and
 Caiming Xiong. A Conversational Paradigm for Program Synthesis. *arXiv preprint*, 2022.
- 740 OpenAI. GPT-4 Technical Report, 2023.

742

743

- Kexin Pei, Zhou Xuan, Junfeng Yang, Suman Jana, and Baishakhi Ray. Trex: Learning execution semantics from micro-traces for binary similarity. *arXiv preprint arXiv:2012.08680*, 2020.
- Kexin Pei, Jonas Guan, Matthew Broughton, Zhongtian Chen, Songchen Yao, David Williams-King, Vikas Ummadisetty, Junfeng Yang, Baishakhi Ray, and Suman Jana. StateFormer: Fine-grained type recovery from binaries using generative state modeling. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 690–702, 2021.
- Alec Radford and Karthik Narasimhan. Improving Language Understanding by Generative Pre Training. 2018. URL https://api.semanticscholar.org/CorpusID:49313245.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
 Models are Unsupervised Multitask Learners. 2019.
- 755 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton,

756 Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel 758 Synnaeve. Code Llama: Open Foundation Models for Code, 2023. 759 Andreas Sæbjørnsen, Jeremiah Willcock, Thomas Panas, Daniel Quinlan, and Zhendong Su. Detect-760 ing Code Clones in Binary Executables. In Proceedings of the Eighteenth International Symposium 761 on Software Testing and Analysis, ISSTA '09, pp. 117-128, New York, NY, USA, 2009. Associa-762 tion for Computing Machinery. ISBN 9781605583389. doi: 10.1145/1572272.1572287. URL 763 https://doi.org/10.1145/1572272.1572287. 764 Calif. Sausalito. Cybercrime Bytes: 10 Hot Security Certs, Public Safety Hacked, Intru-765 sion's Shield. Cybercrime Magazine, 2020. URL https://cybersecurityventures.com/ 766 cybercrime-bytes-10-hot-security-certs-public-safety-hacked-intrusions-shield/. 767 768 Ensheng Shi, Yanlin Wang, Lun Du, Hongyu Zhang, Shi Han, Dongmei Zhang, and Hongbin Sun. 769 CAST: Enhancing Code Summarization with Hierarchical Splitting and Reconstruction of Abstract 770 Syntax Trees. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau 771 Yih (eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 4053–4062, Online and Punta Cana, Dominican Republic, November 2021. 772 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.332. URL https: 773 //aclanthology.org/2021.emnlp-main.332. 774 775 Chad Spensky, Hongyi Hu, and Kevin Leach. LO-PHI: Low-Observable Physical Host Instrumenta-776 tion for Malware Analysis. 2016. 777 Benjamin Steenhoek, Md Mahbubur Rahman, Richard Jiles, and Wei Le. An Empirical Study of Deep 778 Learning Models for Vulnerability Detection. In 2023 IEEE/ACM 45th International Conference 779 on Software Engineering (ICSE), pp. 2237–2248, 2023. doi: 10.1109/ICSE48619.2023.00188. 781 Zian Su, Xiangzhe Xu, Ziyang Huang, Zhuo Zhang, Yapeng Ye, Jianjun Huang, and Xiangyu Zhang. 782 CodeArt: Better Code Models by Attention Regularization When Symbols Are Lacking, 2024. 783 Hanzhuo Tan, Qi Luo, Jing Li, and Yuqun Zhang. LLM4Decompile: Decompiling Binary Code with 784 Large Language Models, 2024. 785 786 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 787 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language 788 Models. arXiv preprint arXiv:2302.13971, 2023. 789 790 Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. Journal of Machine Learning Research, 9(86):2579-2605, 2008. URL http://jmlr.org/papers/v9/ 791 792 vandermaaten08a.html. 793 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N 794 Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett 796 (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Asso-797 ciates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/ 798 3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf. 799 Hao Wang, Wenjie Qu, Gilad Katz, Wenyu Zhu, Zeyu Gao, Han Qiu, Jianwei Zhuge, and Chao Zhang. 800 JTrans: Jump-Aware Transformer for Binary Code Similarity Detection. In Proceedings of the 801 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2022, pp. 802 1-13, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393799. 803 doi: 10.1145/3533767.3534367. URL https://doi.org/10.1145/3533767.3534367. 804 805 Yi Wu, Nan Jiang, Hung Viet Pham, Thibaud Lutellier, Jordan Davis, Lin Tan, Petr Babkin, and 806 Sameena Shah. How Effective Are Neural Networks for Fixing Security Vulnerabilities. In 807 Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2023, pp. 1282–1294, New York, NY, USA, 2023. Association for Computing Machinery. 808 ISBN 9798400702211. doi: 10.1145/3597926.3598135. URL https://doi.org/10.1145/ 809

3597926.3598135.

810 Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. Automated Program Repair in the Era 811 of Large Pre-Trained Language Models. In Proceedings of the 45th International Conference on 812 Software Engineering, ICSE '23, pp. 1482–1494. IEEE Press, 2023. ISBN 9781665457019. doi: 813 10.1109/ICSE48619.2023.00129. URL https://doi.org/10.1109/ICSE48619.2023.00129. 814

- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient Streaming 815 Language Models with Attention Sinks. In *The Twelfth International Conference on Learning* 816 *Representations*, 2024. URL https://openreview.net/forum?id=NG7sS51zVF. 817
 - Xiangzhe Xu, Shiwei Feng, Yapeng Ye, Guangyu Shen, Zian Su, Siyuan Cheng, Guanhong Tao, Qingkai Shi, Zhuo Zhang, and Xiangyu Zhang. Improving Binary Code Similarity Transformer Models by Semantics-Driven Instruction Deemphasis. In Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2023, pp. 1106–1118, New York, NY, USA, 2023a. Association for Computing Machinery. ISBN 9798400702211. doi: 10.1145/3597926.3598121. URL https://doi.org/10.1145/3597926.3598121.
- 825 Xiangzhe Xu, Zhuo Zhang, Shiwei Feng, Yapeng Ye, Zian Su, Nan Jiang, Siyuan Cheng, Lin Tan, 826 and Xiangyu Zhang. LmPa: Improving Decompilation by Synergy of Large Language Model and Program Analysis, 2023b.
- 829 Zhaoyan Xu, Antonio Nappa, Robert Baykov, Guangliang Yang, Juan Caballero, and Guofei Gu. AUTOPROBE: Towards Automatic Active Malicious Server Probing Using Dynamic Binary 830 Analysis. In Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, CCS '14, pp. 179–190, New York, NY, USA, 2014. Association for Computing Machinery. 832 ISBN 9781450329576. doi: 10.1145/2660267.2660352. URL https://doi.org/10.1145/ 833 2660267.2660352. 834
 - Zichao Yang, Divi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical Attention Networks for Document Classification. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1480–1489, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/ N16-1174. URL https://aclanthology.org/N16-1174.
 - Wang Yue, Wang Weishi, Joty Shafiq, and C.H. Hoi Steven. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, 2021.
 - Zhuo Zhang, Yapeng Ye, Wei You, Guanhong Tao, Wen-chuan Lee, Yonghwi Kwon, Yousra Aafer, and Xiangyu Zhang. OSPREY: Recovery of Variable and Data Structure via Probabilistic Analysis for Stripped Binary. In 2021 IEEE Symposium on Security and Privacy (SP), pp. 813–832, 2021. doi: 10.1109/SP40001.2021.00051.
- 849 850 851

818

819

820

821

822

823

824

827

828

831

835

836

837

838

839

840

841

842

843

844 845

846

847

848

APPENDIX Α

852 853 854

A.1 DATA COLLECTION

855 This section provides additional details of the data collection. To collect assemblies from The-Stack, 856 we attempt to compile 4 million C programs, of which 138.8K is compiled successfully. We do not 857 collect more due to the computation resource limitations. 858

For the 757.1K and 138.8K source code that successfully compiled into executables (using gcc) 859 from AnghaBench and The-Stack, we disassemble them using objdump. objdump was not able to 860 successfully disassemble all the executables, resulting in some empty assembly code. Thus, the 861 number of 00 - 01 we obtain from each corpus is different and smaller than the number of source 862 codes as shown in Table 1. 863

Figure 6 shows an example of preprocessing the raw assembly code as described in Section 3.1.



Figure 6: Example of assembly code preprocessing

Table 10: Comparison with applying hierarchical attention on all attention heads, using 1B models.

			Pass@1					Pass@10		
Techniques	00	01	O2	O3	Avg.	00	01	O2	O3	Avg.
Nova-HA	30.58	19.88	20.58	16.40	21.86	44.75	33.13	33.31	29.82	35.25
Nova (hierarchical on all heads)	32.38	18.87	20.56	16.34	22.04	45.95	32.19	32.78	29.01	34.98
Nova	37.53	21.71	22.68	18.75	25.17	49.38	34.84	36.95	32.03	38.30

A.2 TRAINING DETAILS

This section provides additional details of training. We pre-train Nova starting from DeepSeek-Coder, using the language modeling objective (L_{lm}) for one epoch on the C functions and assembly functions collected from AnghaBench and The-Stack corpora. This step uses a batch size of 128, with the input truncated by a 1,024 tokens limit. The model weights are updated using the AdamW optimizer. The learning rate is $5e^{-5}$, using 1000 steps of warm-up and a cosine decay to adjust the learning rate.

Then, the model is further pre-trained with the combination of language modeling and contrastive learning objectives $(L = L_{lm} + \lambda(L_{fcl} + L_{ocl}))$, with λ set to 0.1. To train with the functionality contrastive learning objective, we discard any source code that misses any one of 00 – 03 assemblies and also discard the source code whose 02 assembly is the same as its 03 assembly. As a result, this step is only trained for 0.36M data samples for one epoch. The batch size is 64, with the input truncated by a 1,024 tokens limit. The learning rate is $2e^{-5}$ using the AdamW optimizer.

The fine-tuning of both BCD and BCSD uses a batch size of 64, with the input truncated by a 2,048 token limit. Similarly, the learning rate is $2e^{-5}$ using the AdamW optimizer, and the model is fine-tuned for one epoch. During the training using the contrastive learning objectives, and the fine-tuning of BCSD, we use the average of [INST] tokens' last layer hidden states to represent the embedding of a binary function.

Infrastructure The training are conducted on eight NVIDIA RTX A100 GPUs, each with 40GB memory. Our implementation is based on Huggingface's implementation of DeepSeek-Coder ¹, PyTorch ², and DeepSpeed ³.

903 904

905

899

872

873 874

882

A.3 APPLYING HIERARCHICAL ATTENTION ON HALF ATTENTION HEADS

The hierarchical attention mask is always applied on half of the attention heads at each layer in
 Nova. This ensures the LLM balances the hierarchical knowledge of assembly code and pre-trained
 knowledge learned by full self-attention.

We conducted experiments applying hierarchical attention to all the attention heads. Results in
Table 10 show that when applying hierarchical attention to all the attention heads of transformer
layers, the performance does not improve enough and even drops under some settings. This implies
that the standard full self-attention mechanism indeed learns knowledge that may not be captured by
hierarchical attention. Thus, we only apply hierarchical attention to half of the attention heads in each
transformer layer of Nova to balance the knowledge learned by standard and hierarchical attention.

915 916

917

³https://github.com/microsoft/DeepSpeed

¹https://huggingface.co/deepseek-ai/deepseek-coder-1.3b-base

²https://pytorch.org/get-started/locally/

A.4 ADDITIONAL ABLATION STUDY ON BINARY CODE DECOMPILATION

We provide additional ablation studies on studying the impact of each individual contrastive learning objective. We study two more models:

Nova_FCL-HA: Removing functional contrastive learning and hierarchical self-attention.

Nova_OCL-HA: Removing optimization contrastive learning and hierarchical self-attention.

Table 11 shows that each component of the contrastive learning brings certain improvements to the Pass@1 and Pass@10 on HumanEval-Decompile, and we find the impact of function contrastive learning (FCL) is larger than the impact of optimization contrastive learning (OCL), suggesting that aligning the model's embeddings for assembly code with the same functionality is more useful.

Table 11: Ablation study of Nova-1B on HumanEval-Decompile.

			Pass@1					Pass@10		
Techniques	00	01	02	O3	Avg.	00	01	02	O3	Avg.
Nova_CL-HA	20.73	16.16	15.03	11.19	15.78	33.55	28.12	26.96	21.01	27.41
Nova-FCL-HA	22.38	16.20	16.37	13.25	17.05	36.13	29.48	30.02	23.76	29.85
Nova_OCL-HA	28.44	18.87	18.53	15.76	20.40	40.28	32.33	31.80	27.05	32.87
$Nova_{-HA}$	30.58	19.88	20.58	16.40	21.86	44.75	33.13	33.31	29.82	35.25
Nova	37.53	21.71	22.68	18.75	25.17	49.38	34.84	36.95	32.03	38.30

<func< th=""><th>0>:</th><th></th><th><pre>int func0(float *arr, int n) {</pre></th></func<>	0>:		<pre>int func0(float *arr, int n) {</pre>
<func 0: 4: 6: 8: b: 14: 15: 14: 16: 17: 27: 24: 22a: 22a: 32: 32: 33: 33:</func 	0>: endbr64 test jle lea add mov movss 00 cmp je movss mov movaps subss andps comiss ja add	<pre>4 %esi, %esi 48 <func0+0x48> -0x1(%rsi), %ecx \$0x1, %rcx \$0x1, %edx 0x0(%rip), %xmm3 %rcx, %rdx 4e <func0+0x4e> -0x4(%rdi, %rdx, 4), %xmm1 %xmm1, %xmm1 %xmm1, %xmm1 %xmm1, %xmm0 54 <func0+0x54> \$0x1 % rex</func0+0x54></func0+0x4e></func0+0x48></pre>	<pre>int func@(float *arr, int n) { for (int i = 1; i < n; i++) { for (int j = 0; j < i; j++) { if (fabs(arr[i] - arr[j]) < 1.0f return 1; } } return 0; GPT-4's Incorrect Decompilation int func@(const float *x, int n) { int i; float diff = 1.0f; for (i = 1; i < n; i++) { if (fabs(x[i] - x[i-1]) > diff) return 0; return 1; } } </pre>
3a: 3e: 40.	add cmp io	<pre>\$0X1, %Fax %eax, %esi 2a <func0+0x2a></func0+0x2a></pre>	LLM4Decompile's Incorrect Decompilation
42: 46: 48: 4d: 4e: 53: 54: 59:	add jmp mov retq mov retq mov retq retq	<pre>\$0x1, %rdx 1c <func0+0x1c> \$0x0, %eax \$0x0, %eax \$0x1, %eax</func0+0x1c></pre>	<pre>int func@(float *a, int n, float e) { int i, j; for (i = 0; i < n; i++) for (j = i + 1; j < n; j++) if (fabs(a[i] - a[j]) < e) return 1; return 0; }</pre>
		O1 Assembly Code	Peturn 0; Nova's Correct } Decompilation

Figure 7: Nova-1B correctly decompiles HumanEval-Decompile task 0.

A.5 BINARY CODE DECOMPILATION CASE STUDIES

Figure 7 shows an example from HumanEval-Decompile (task_id 0). Given the 01-optimized assembly code, GPT-4 fails to figure out the number of function arguments correctly, missing one important argument "float e", and thus produces wrong functionality in the decompiled code. LLM4Decompile-1B makes similar mistakes and also misses the inner nested for loop. Nova-1B correctly decompiles the assembly into source code, where the ground truth is checking if any two elements in the given list \star a (with size n) are close to each other than a given threshold e.

Figure 8 shows another more complex example, HumanEval-Decompile task_id 20. Nova-1B correctly decompiles the source code, successfully figuring that the function is trying to find the two elements that are closest to each other in the given array *a, with the result stored in minmax.

972			
973			
974			
975	<func0>:</func0>		<pre>void func0(float *a, int n, float *b) {</pre>
	0: endbr64	1	int i, j;
976	4: movss	(%rdi), %xmm0	float min, max;
977	8: movss	%xmm0, (%rdx) 0x4(%rdi) %xmm1	$\min - \max - 2[0]$
978	11: movss	%xmm1, 0x4(%rdx)	for $(i = 1; i < n; i++)$ {
070	16: test	%esi, %esi	if (a[i] < min)
979	18: jle	8b <func0+0x8b></func0+0x8b>	<pre>min = a[i];</pre>
980	1a: lea	-0x1(%rsi), %r9d	<pre>else if (a[i] > max)</pre>
	le: mov	\$0x1, %r8d	<pre>max = a[1];</pre>
981	24: mov	%rdi, %rcx	}
982	2/: movss	0x0(%rip), %xmm2	b[0] = min;
	,2e. 00	¢01 %0	LLM4Decompile's Incorrect
983	33: movss	\$0X1, %Γ9 0x0(%rin) %xmm3	³ Decompilation
984	3a: 00		
985	3b: cmp	%r8, %r9	<pre>void func0(float *a, int n, float</pre>
505	;3e: je	82 <func0+0x82></func0+0x82>	*minmax) {
986	40: mov	%r8, %rax	minmax[0] = a[0];
097	43: nopl	0x0(%rax, %rax, 1)	minmax[1] = a[1];
907	48: movss	(%rcx), %xmm1	
988	4c: movaps	%xmm1, %xmm0	for (int $1 = 0; 1 < n; 1++)$ {
080	4I: SUDSS	(%rdl, %rax, 4), %xmm0	float diff = fabs(aFi] = aFi];
909	57: comiss	%×mm0, %×mm2	if (diff < minmax[1] - minmax[0])
990	5a: jbe	6d <func0+0x6d></func0+0x6d>	{
991	5c: movss	%xmm1, (%rdx)	minmax[0] = a[i];
002	60: movss	(%rdi, %rax, 4), %xmm1	<pre>minmax[1] = a[j];</pre>
552	68: moves	%xmm1 Qx4(%ndx)	
993	6d add	\$0x1 $$rax$	}
994	71: cmp	%eax, %esi	
004	73: jg	48 <func0+0x48></func0+0x48>	if (minmax[0] > minmax[1]) {
995	75: add	\$0×1, %r8	<pre>float temp = minmax[0];</pre>
996	79: add	\$0x4, %rcx	<pre>minmax[0] = minmax[1]; minmax[1] = temp;</pre>
997	99: retq		} Nova's Correct
998	02	Assembly Code	} Decompilation

Figure 8: Nova-1B correctly decompiles HumanEval-Decompile task 20.

Benchm	arks N	ova_ <i>CL</i> - <i>HA</i>	Nova _{-HA}	Nova-1B
Binutils		0.86	0.88	0.87
Curl		0.84	0.87	0.89
ImageM	agick	0.79	0.80	0.86
SQLite		0.80	0.83	0.77
OpenSS	L	0.90	0.92	0.90
Putty		0.68	0.66	0.72
Avg.		0.81	0.83	<u>0.84</u>

Benchmarks	Nova $_{-CL-HA}$	Nova $_{HA}$	Nova-1B
Binutils	0.80	0.82	0.79
Curl	0.84	0.84	0.86
ImageMagick	0.70	0.72	0.79
SQLite	0.74	0.78	0.73
OpenSSL	0.89	0.89	0.88
Putty	0.59	0.60	0.65
Avg.	0.76	0.78	0.78

Table 12: Ablation study with K = 50.

Table 13:	Ablation	study	with	K =	= 100.
10010 15.	1 ionation	Study	** 1111	17 -	- 100.

Benchmarks	Nova_CL-HA	Nova $_{HA}$	Nova-1B
Binutils	0.71	0.74	0.73
Curl	0.80	0.73	0.83
ImageMagick	0.61	0.63	0.73
SQLite	0.68	0.71	0.68
OpenSSL	0.85	0.87	0.84
Putty	0.53	0.53	0.55
Avg.	0.70	0.70	<u>0.73</u>

Table 14: Ablation study with K = 200.

Nova_HA Nova-1B Benchmarks Nova $_{-CL-HA}$ Binutils 0.62 0.65 0.65 Curl 0.67 0.71 0.73 ImageMagick 0.46 0.51 0.61 0.61 0.59 SQLite 0.62 OpenSSL 0.77 <u>0.79</u> 0.78Putty 0.46 0.46 0.47

Table 15: Ablation study with K = 500.

0.62

<u>0.64</u>

0.60

1024 1025

999

1012 1013

1014 1015

1016

1017

1018

1019

1020

1021

1022

1023

Avg.

1026 A.6 BINARY CODE SIMILARITY DETECTION ABLATION STUDY

Table 12, 13, 14, 15 show the detailed ablation study results of BCSD. Nova wins on the most benchmarks when K = 100 or 500, and ties with Nova_{-HA} when K = 50, or 200.



Figure 9: PCA of embeddings calculated by Nova-CL-HA, Nova-HA, and Nova.

1043 A.7 Additional Analysis of Embedding

Figure 9 shows the results of PCA of embeddings provided by $Nova_{-CL-HA}$, $Nova_{-HA}$, and Nova, on randomly sampled seven examples, where Nova's embeddings are consistently more distinguishable by functionalities.







1080 A.8 Additional Analysis of Attention