

# Multi-modal Learning for WebAssembly Reverse Engineering

Hanxian Huang

University of California at San Diego  
San Diego, USA  
hah008@ucsd.edu

Jishen Zhao

University of California at San Diego  
San Diego, USA  
jzhao@ucsd.edu

## ABSTRACT

The increasing adoption of WebAssembly (Wasm) for performance-critical and security-sensitive tasks drives the demand for WebAssembly program comprehension and reverse engineering. Recent studies have introduced machine learning (ML)-based WebAssembly reverse engineering tools. Yet, the generalization of task-specific ML solutions remains challenging, because their effectiveness hinges on the availability of an ample supply of high-quality task-specific labeled data. Moreover, previous works trained models only with features extracted from WebAssembly, overlooking the high-level semantics present in the corresponding source code and its documentation. Acknowledging the abundance of available source code with documentation, which can be compiled into WebAssembly, we propose to learn representations of them concurrently and harness their mutual relationships for effective WebAssembly reverse engineering.

In this paper, we present WasmRev, the first multi-modal pre-trained language model for WebAssembly reverse engineering. WasmRev is pre-trained using self-supervised learning on a large-scale multi-modal corpus encompassing source code, code documentation and the compiled WebAssembly, without requiring labeled data. WasmRev incorporates three tailored multi-modal pre-training tasks to capture various characteristics of WebAssembly and cross-modal relationships. WasmRev is only trained once to produce general-purpose representations that can broadly support WebAssembly reverse engineering tasks through few-shot fine-tuning with much less labeled data, improving data efficiency. We fine-tune WasmRev onto three important reverse engineering tasks: type recovery, function purpose identification and WebAssembly summarization. Our results show that WasmRev pre-trained on the corpus of multi-modal samples establishes a robust foundation for these tasks, achieving high task accuracy and outperforming the state-of-the-art ML methods for WebAssembly reverse engineering.

## CCS CONCEPTS

- **Software and its engineering** → **Software notations and tools**; • **Security and privacy** → **Software reverse engineering**;
- **Computing methodologies** → **Machine learning**.

## KEYWORDS

WebAssembly, reverse engineering, representation learning, multi-modal learning, program language modeling, function purpose identification, type recovery, code summarization.

### ACM Reference Format:

Hanxian Huang and Jishen Zhao. 2024. Multi-modal Learning for WebAssembly Reverse Engineering. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA '24)*, September 16–20, 2024, Vienna, Austria. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3650212.3652141>

## 1 INTRODUCTION

WebAssembly (Wasm) is a low-level, portable, bytecode format compiled from high-level languages, such as C, C++, and Rust, delivering near-native performance when executed on the web [1, 29]. It is a promising technology to enhance performance of various applications traditionally developed in JavaScript, e.g., cryptography [6, 71] and machine learning [18, 20]. Initially conceived for execution in web browsers, WebAssembly is now expanding its reach to encompass various application domains, such as serverless computing [3, 24], edge computing [23], and IoT [48, 72].

As WebAssembly gains popularity across diverse applications, an increasing demand emerges for understanding and reverse engineering WebAssembly code. Many WebAssembly modules – including potentially malicious ones – are distributed through third-party services, rendering the source code unavailable on the client-side [51]. This requires users to understand the WebAssembly modules and audit it to prevent potential attacks or malicious code [40, 62, 70]. While WebAssembly has a readable text format, manual interpretation and comprehension remain challenging and error-prone. Different from register-based native binaries (e.g., x86 or ARM), WebAssembly adopts a stack machine design, necessitating the tracing of stack behavior to comprehend the code or compute a specific variable. Moreover, WebAssembly only employs four numeric data types, i32, i64, f32, and f64, which conceals the data type information and further increases the complexity of comprehending WebAssembly.

Several prior studies explored WebAssembly understanding and reverse engineering with conventional analysis approaches [8, 22, 34, 65] or machine learning (ML)-based methods [42, 61]. Traditional approaches analyze WebAssembly through precise, logical reasoning, often incorporating heuristics to ensure the practical utility of the tools [8, 22, 34, 65]. However, crafting effective heuristics is difficult, especially in cases where an accurate analysis result depends on uncertain information, such as natural language (NL) content in code documentation and NL code search tasks [21, 31, 47]. This uncertainty is not conducive to logic-based reasoning [60].

Instead, recent ML-based solutions learn features from extensive labeled code datasets and manipulate fuzzy information in code

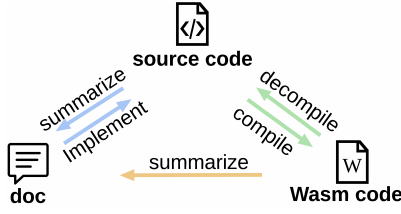
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

ISSTA '24, September 16–20, 2024, Vienna, Austria

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0612-7/24/09

<https://doi.org/10.1145/3650212.3652141>

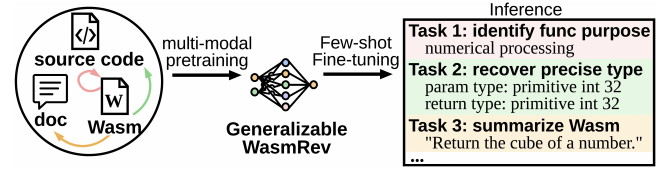


**Figure 1: The relationship of engineering and reverse engineering among source code, code documentation, and WebAssembly code.**

without relying on manually encoded precise rules. SnowWhite [42] trains a sequence model based on WebAssembly code to recover precise types. WASPur [61] is trained on control flow graph features from WebAssembly code to identify function purpose. However, one key caveat is that these frameworks are designed for specific tasks, while overlooking the generic WebAssembly features shared across various WebAssembly analysis tasks. In addition, both frameworks require large scale task-specific labeled datasets, which are laborious to collect for many analysis tasks. For instance, function purpose identification requires manual inspection of WebAssembly code, checking high-level semantics like function names to label the correct purpose of a function [61]. Hence, generalizing and applying these task-specific models to a broader range of WebAssembly tasks poses a significant challenge.

Despite substantial efforts in training ML models for WebAssembly [42, 61], the relationships among high-level source code, code documentation, and WebAssembly code remain crucial yet under-explored in the realm of WebAssembly reverse engineering. As depicted in Figure 1, a program is often composed of both code snippets (in a high-level programming language (PL)) and corresponding documentation (in natural language (NL)), while the code snippets can be compiled into WebAssembly (Wasm). In this paper, we define source code, code documentation, and WebAssembly code as multiple modalities of code – (NL, PL, Wasm). The source code (PL) serves as a human-readable representation of the intended logic of a program in a high-level programming language; the code documentation (NL) supplements this by offering insights into the purpose of functions in natural language. Studying both offers contextual information that aids in interpreting the purpose and behavior of the code, providing valuable knowledge for reverse engineers seeking a comprehensive understanding and interpretation of WebAssembly code. As such, concurrently learning representations of these semantically equivalent modalities, along with understanding their relationships, will yield complementary information critical for comprehending WebAssembly.

To address the aforementioned challenges and limitations, we propose WasmRev, the first multi-modal pre-trained language model for WebAssembly reverse engineering. Figure 2 shows an example of WasmRev workflow. Different from previous task-specific models [42, 61], WasmRev effectively learns a generic representation among WebAssembly code, source code, and code documentation, and efficiently transfers it to perform various WebAssembly reverse engineering tasks. To this end, we tailor three pre-training tasks enabling WasmRev to effectively learn inter-modal and intra-modal relationships, and close the gaps across various modalities of representations (Section 4.3). The pre-training is self-supervised, without



**Figure 2: WasmRev overview.**

requiring a large scale labeled dataset, effectively solving the inherent scarcity of labeled dataset issues. The pre-training is done once, while the pre-trained WasmRev can be shared and adapted to various WebAssembly tasks with light-weight fine-tuning. Compared with supervised learning on specific tasks, such a process is more efficient and requires much less task-specific labeled data (Section 5). WasmRev offers efficient and accurate solutions for various WebAssembly reverse engineering tasks, furnishing programmers with insightful high-level abstractions to enhance an accessible comprehension of WebAssembly and facilitate further code inspection. In summary, this paper contributes the following:

- We propose WasmRev, a pre-trained language model that learns a generic representation from multi-modal inputs, including WebAssembly code, source code, and documentation. This generalizable representation can facilitate understanding of low-level languages like WebAssembly and can be effectively transferred to diverse WebAssembly reverse engineering tasks. To the best of our knowledge, this is the first multi-modal pre-trained language model for generalized WebAssembly reverse engineering.
- Enabled by WasmRev, we develop an ML-based WebAssembly type recovery tool and a function purpose identification tool, which demand much less labeled data. In addition, we develop the first WebAssembly summarization tool. Together, these tools deliver accurate analysis results, empowering programmers with insights needed for thorough WebAssembly comprehension and further code inspection.
- We perform a comprehensive experimental evaluation on our WebAssembly reverse engineering tools. The experiment results reveal that WasmRev achieves high accuracy and data efficiency across all tasks, surpassing state-of-the-art (SOTA) ML methods for WebAssembly.

## 2 BACKGROUND AND MOTIVATION

We use an motivating example to discuss the challenges in WebAssembly understanding, the limitations of state-of-the-art methods, and how our design tackles these challenges and limitations.

### 2.1 WebAssembly Basics

The WebAssembly standard [1] defines assembly-like bytecode with a unique instruction set architecture (ISA) and specific binary encodings for various types of operations. A WebAssembly binary is a binary encoding of a module with a clear modular structure composed of several sections, such as parameters and result types, functions, and data. To ensure readability, the standard provides a text format that offers a readable representation of the internal structure of a module, including type, memory, and function definitions. WebAssembly programs are commonly compiled from high-level programming languages (e.g., C, C++, and Rust) using

WebAssembly compilers, such as Emscripten [15] and Rustc [63]. Figure 3 (a) and (b) show examples of a C source code and its corresponding compiled WebAssembly in text format.

## 2.2 Challenges in WebAssembly Understanding

While WebAssembly offers promising performance and portability, the ecosystem of its reverse engineering tools is still in an early stages of development. Substantial challenges remain for WebAssembly programmers, particularly those who are novices.

**Motivating example.** When the high-level source code is accessible (Figure 3(a)), a programmer can comprehend the corresponding WebAssembly code (Figure 3(b)) by cross-referencing it with the relevant high-level code snippet and documentation; debugging tools such as DWARF [14] may also facilitate the comprehension. However, a common scenario often involves a WebAssembly binary delivered from a third-party, with no access to the high-level source code, source map, or debugging information. As a result, programmers will need to be familiar with the WebAssembly specifications and manually deduce the code behavior; alternatively, programmers may seek assistance from reverse engineering tools, such as WABT [26], to convert the binary into more readable high-level code (Figure 3(c)) or decompile it into high-level syntax (Figure 3(d)).

We identify two major challenges with the current approaches to understanding WebAssembly. First, **(C1) a lack of high-level information** (Figure 3(b)(c)(d)). High-level programming languages (Figure 3(a)) carry informative descriptions and convey a program's intent via function names, variable names, and natural language comments and documentation. But WebAssembly code purely consists of low-level instructions, operating at a low level of abstraction and depicting stacking behavior. Moreover, WebAssembly only employs four data types, i32, i64, f32, and f64, significantly concealing the data type information. For instance, an i32 type may represent a signed or unsigned integer, an array, a pointer to a struct, or many other source types. In our example, neither the converted C source (Figure 3(c)) nor the decompiled code (Figure 3(d)) successfully reconstructs lexical information or precise types: Figure 3(c) only recovers u32 or u64, which is uninformative; Figure 3(d) recovers the input parameter type as int, which is ambiguous because the parameter is sometimes treated as a pointer or an array within the function body, e.g., a[0]. These gaps in representation present hurdles to precise code reasoning and inspection. Second, **(C2) tracking stack behavior is cumbersome and error-prone.** WebAssembly execution is based on a stack machine architecture, which requires carefully deductive reasoning to track stack behavior. The absence of detailed debugging information and the presence of unknown optimizations substantially complicate the task of comprehending and interpreting WebAssembly. This is also less intuitive for developers accustomed to the widely adopted register-based native architectures. In our example, the converted C source (Figure 3(c)) introduces intermediate variables and simulates the stack behavior with them using the high-level C code. However, this still requires programmers to trace the stack behavior by monitoring these variables, a process that is both tedious and error-prone, especially for non-experts. Furthermore, the complexity of stack tracing increases with program complexity, rendering manual inspection of real-world projects increasingly impractical.

To overcome these challenges, it is imperative to furnish programmers with higher-level information, such as high-level types, function purposes, and code summarization. Such information is particularly crucial for inspecting stripped or minified WebAssembly modules that lack high-level semantics. This higher-level information, coupled with the decompiled high-level syntax or the converted high-level source code, will offer more comprehensive insights for reverse engineering, help programmers understand WebAssembly more effectively, and reduce the steep learning curve for novices.

## 2.3 Limitations in the State-of-the-Art Methods

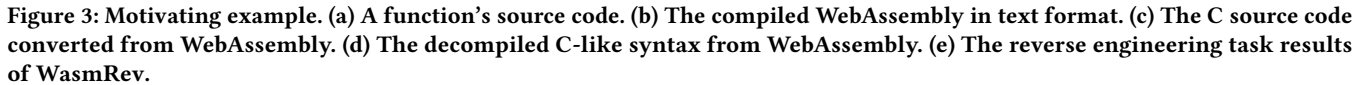
**Conventional reverse engineering.** Conventional methods [8, 22, 26, 34, 65] employ logical reasoning and rely on heuristics for WebAssembly analysis. The heuristics are hand-coded by experts, requiring non-trivial development effort to create, maintain, and update, especially considering the continuous evolution of WebAssembly and compilers. Furthermore, logic-based methods are unsuitable for tasks such as code summarization, which involves uncertain information (e.g., natural language).

**ML-based WebAssembly reverse engineering.** Machine learning solutions, on the contrary, automatically learn heuristics and underlying features from large datasets of code and are good at handling fuzzy information [60]. Given the promising capabilities of ML models, researchers have explored task-specific, supervised learning solutions for WebAssembly [42, 61], involving collecting a task-specific labeled dataset and supervised training a task-specific model. Despite the promising task performance, these studies impose several limitations. **(L1) Requiring large-scale labeled datasets:** Most resources are unlabeled. Despite the increasing number of open-source repositories in collaborative developer platforms like GitHub [2, 49], most of the collected WebAssembly snippets are unlabeled, which can not be directly used by current supervised ML approaches. The inherent scarcity of labeled data for many WebAssembly analysis tasks necessitates extensive manual efforts for data collection, e.g., manually inspecting and labeling the purpose of functions [61]. **(L2) Lacking generalization or transferability:** task-specific models necessitating not only task-specific labeled dataset, but also customizing, training and maintenance of each model for each specific task, which is inefficient and makes them difficult to generalize or transfer to other WebAssembly tasks. **(L3) Overlooking available high-level semantics:** Existing ML methods only utilize high-level information such as high-level types and function purpose as labels in supervised learning. However, they are overlooking additional potential information and context in the high-level semantics. For example, function names and comments could provide clues about the function's intent, while the way a data object is being used in context could hint at the variable type.

## 2.4 Our Design: WasmRev

To address **(L1)**, inspired by the success of self-supervised pre-training methods in natural language processing (NLP) and software engineering [19, 21, 38, 69], we train WasmRev to learn a generalizable WebAssembly representation through self-supervised learning. The essence of self-supervised pre-training lies in generating “pseudo-labels” without the need for manual labeling. This can be easily achieved by designing “pseudo-tasks”, for example,





To tackle (C1) and (L3), we propose to learn representations of code documentation, source code and WebAssembly code concurrently – (NL, PL, Wasm) multi-modal learning. Specifically, we design a pre-training task that guides the model to predict missing words (which can be in NL, PL, Wasm) in randomly masked multi-modal inputs, facilitating the model in learning cross-modal relationships and dependencies (Section 4.3 T1). We design another task to identify similar semantics, i.e., various modalities of the same program while distinguishing different semantics, so as to bridge the gaps among different modal representations. Developing such multi-model representation learning is non-trivial, due to the different language specifications and grammars. WebAssembly employs rigid instructions to represent stack behaviour,

with instruction order significantly influenced by the compiler and optimization levels, thereby increasing contextual complexity. In contrast, source code employs high-level syntax and follows specific coding conventions. Natural language is weakly structured and varies in style among programmers or documentation creators. These factors collectively contribute to the complexity of processing multiple languages within a single neural architecture. Despite the success of language modeling in source code languages [21, 28] and native binaries [43, 69], given the language differences, previous solutions are not ideally suitable or directly applicable to WebAssembly. To the best of our knowledge, no prior attempt has been made to learn a multi-modal representation of (NL, PL, Wasm) for WebAssembly reverse engineering.

To address the challenges and limitations, we propose WasmRev, the first multi-modal language model for WebAssembly understanding and reverse engineering. The goal of WasmRev is to learn a generalizable representation, i.e., learn to project multi-modal inputs into a representative embedding space, which can be used effectively in various WebAssembly reverse engineering tasks.

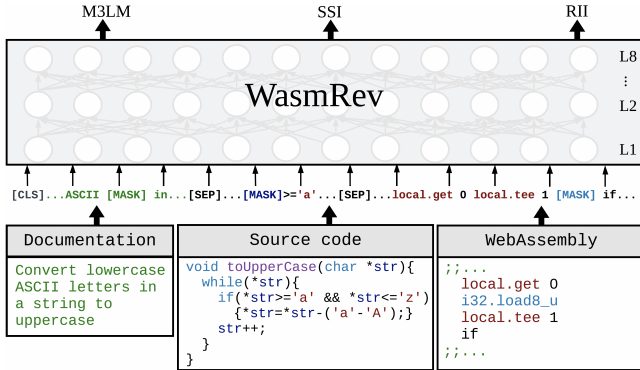
- (1) Multi-modal pre-training (Section 4.3): WasmRev is trained with tailored multi-modal pre-training tasks on our collected (code documentation, source code, WebAssembly code) samples, to learn a generic multi-modal embedding.
- (2) Few-shot fine-tuning (Section 5): The pre-trained WasmRev is allowed for re-purposing it to various WebAssembly reverse engineering tasks. To evaluate the generalization of the pre-trained

WasmRev, we fine-tune it on three tasks: type recovery, function purpose identification and WebAssembly summarization. The fine-tuning is few-shot, requiring much less task-specific labeled datasets than non pre-train methods.

- (3) Inference: Each of the fine-tuned models performs inference on new WebAssembly input, and generates accurate prediction. The collected outputs together serve programmers with high-quality insights who decide to perform further code inspection. An example report is shown in Figure 3 (e).

## 4 LEARN A GENERIC MULTI-MODAL REPRESENTATION

In this section, we will elaborate how WasmRev learns a generic (NL, PL, Wasm) multi-modal representation. As illustrated in Figure 4, WasmRev comprises three key components: WasmRev model, input/output representations and pre-training tasks.



**Figure 4: An illustration of WasmRev pre-training. The model architecture consists of 8 bidirectional Transformer layers. The model takes the concatenation of (documentation, source code, WebAssembly code) as input, and is pre-trained on M3LM, SSI and RII tasks.**

### 4.1 WasmRev Model Architecture

We adopt BERT [19] as WasmRev’s backbone, inspired by its promising performance in various code-related tasks [21, 28]. As the model architecture is standard, we will not review the Transformer architecture in detail. As will be elaborated later, we adapt the model through our tailored pre-training process to be WasmRev for better WebAssembly comprehension.

### 4.2 Input / Output Representations

Previous WebAssembly ML models [42, 61] consider only WebAssembly as input. In this paper, we argue that aggregating complementary information from semantically equivalent modalities – source code, code documentation and WebAssembly code – is benefit for WebAssembly comprehension, potentially better supporting diverse downstream tasks and yielding promising results (Section 7.4). Given a NL documentation  $c = \{c_1, c_2, \dots, c_{|c|}\}$ , we consider it as a sequence of words, and split it as WordPiece [73]. Given the corresponding source code in a high level PL  $s = \{s_1, s_2, \dots, s_{|s|}\}$ , and the compiled WebAssembly  $w = \{w_1, w_2, \dots, w_{|w|}\}$ , we regard

them as sequences of tokens. WasmRev takes the concatenation of multiple modalities (NL, PL, Wasm) as model input, i.e.,

$$x = \{[CLS], c, [SEP], s, [SEP], w, [SEP]\} \quad (1)$$

where [CLS] is a special token at the beginning of the input sequence, whose final hidden representation is considered as the aggregated sequence representation for classification or ranking; and [SEP] is a special token to split two kinds of sub-sequences [19]. We normalize instructions to avoid the out-of-vocabulary problem, following previous studies on native binary representation learning [43, 69]. Specifically, string literals are replaced by a special token [STR]. Large constants, empirically those greater than 0xffff, which likely represent memory addresses, are normalized by a special token [ADDR]. Small constants are retained as individual tokens as they may carry valuable information, such as specifying which local variables or function arguments to access. Then we build a vocabulary upon the normalized corpus and tokenized the inputs. We also incorporate position embedding and segmentation embedding into token embedding, and use the mixed vector of them as model input. Specifically, position embedding represents different positions in the input sequence, while segmentation embedding distinguishes documentation, source code and WebAssembly code.

The output of WasmRev includes (1) contextual vector representation of each token, for code documentation, source code, and WebAssembly code; (2) the representation of [CLS], which works as the aggregated sequence representation.

### 4.3 Pre-training Tasks

To best adapt WasmRev for multi-modal learning, particularly for WebAssembly reverse engineering, we go beyond simply adopting BERT’s original pre-training approach. Instead, we design three self-supervised pre-training tasks that are specifically crafted to exploit and learn from our multi-modal input representation, enabling WasmRev to grasp inter-modal and intra-modal relationships.

**(T1) Multi-Modal Masked Language Model (M3LM)** Our approach jointly models NL, PL, and Wasm, providing complementary information contained in multiple modalities. We first extend the NLP task Masked Language Model (MLM) of BERT to multiple modalities. Given a data point of (NL, PL, Wasm) triplet  $\{c, s, w\}$  as input, we randomly select 15% of tokens from the concatenation of (NL, PL, Wasm). We show an M3MLM example in Figure 5. For the selected tokens, we replace 80% of them with [MASK] tokens, 10% with random tokens, and the remaining 10% is unchanged. Formally, the loss function of M3LM is defined as:

$$\mathcal{L}_{M3LM} = - \sum_{i \in M} \log P(t_i | f_M) \quad (2)$$

where  $M$  is the set of indices of the masked or replaced tokens, and  $f_M$  is the corrupted input,  $P(t_i | f_M)$  is the probability for predicting a particular token  $t_i$ , with  $f_M$  as input, calculated by the representation of the token from WasmRev following a softmax function to normalize the output.

By extending conventional MLM to M3LM, the model is guided to learn not only the intra-modal contextual relationship to reconstruct missing words, but also the cross-modal dependencies. In particular, if the WebAssembly context alone proves insufficient for deducing the masked WebAssembly token, the model can draw

upon both documentation and source code to enhance its understanding. This also holds true for the masked token prediction for both documentation and source code.

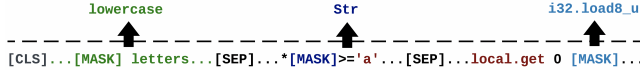


Figure 5: A multi-modal masked language model example

**(T2) Similar Semantics Identification (SSI)** To bridge the gap among different modalities, we task the model with identifying similar semantics, i.e., various modalities of the same program. While the documentation, source code and WebAssembly code of the same program correspond to the same program semantics, they are dissimilar to representations of another program. As such, we create semantically similar and dissimilar pairs as model inputs. There are numerous ways to create these semantically similar pairs, including single-modal pairs like NL vs. Wasm, bi-modal pairs like (NL, PL) vs. (NL, Wasm), and triple-modal pairs like (NL, PL, Wasm) vs. (NL, Wasm, PL). Our selection of semantic-similar pairs is guided by potential downstream applications. Motivated by the WebAssembly summarization task, a Wasm-NL task, we include NL vs. Wasm. Considering the potential WebAssembly decompilation tasks, which are Wasm-PL tasks, we incorporate PL vs. Wasm. With the aim of enriching the variety of similar pair formats to enhance the SSI task, we further include additional bi-modal and triple-modal pairs. We include (NL, PL) vs. (NL, Wasm) pairs, which encourage the model to learn the semantics of PL and Wasm, with NL providing context. Similarly, we include (NL, PL) vs. (PL, Wasm) to enable our model to learn the semantics behind NL and Wasm, with PL as context. We also include (NL, PL, Wasm) vs. (NL, Wasm, PL) to help the model to learn semantics behind different orders of modalities. Furthermore, the samples with WebAssembly code compiled from the same source code but with different compilation flags are also semantically similar pairs.

We then employ in mini-batch and cross mini-batch sampling strategies [10] to construct dissimilar samples. For a batch of training data  $b_1 = [x_1 \dots x_N]$  with batch size  $N$ , we can first obtain a batch of data  $b_2 = [x_1^+ \dots x_N^+]$  where  $(x_i, x_i^+)$  is a pair of similar semantics as describe above, e.g., NL vs. Wasm. The dissimilar semantics  $x_i^-$  of  $x_i$  is chosen as  $x_j$  and  $x_j^+$ ,  $\forall i \neq j$  and  $x_i$  and  $x_j/x_j^+$  are corresponding to different source code. We show an example of SSI in Figure 6. For an input  $x_i$  with representation  $v_i$ , the goal of SSI is to maximize the representation similarity (dot product) between similar samples, while minimizing the representation similarity between dissimilar samples. Formally, the loss of SSI is defined as:

$$\mathcal{L}_{SSI} = - \sum_{i \in \{N, N^+\}} \ln \frac{\exp(v_i \cdot v_i^+)}{\exp(v_i \cdot v_i^+) + \sum_{k=1}^{|x_i^-|} \exp(v_i \cdot v_k^-)} \quad (3)$$

where  $N$  and  $N^+$  are the sets of all program samples covering similar samples discussed above.

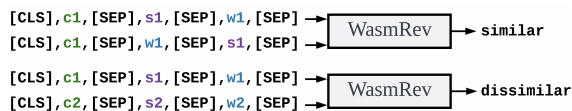


Figure 6: A similar semantics identification examples

**(T3) Reordered Instructions Identification (RII)** Since our target is WebAssembly understanding, we design a task specifically aimed at enabling the model to learn WebAssembly representation, especially the contextual and control flow features. In WebAssembly, the control flow is well-structured. Due to the stack-based architecture of WebAssembly, each instruction controls the next instruction in code. As such, we design a reordered instructions identification task to help WasmRev comprehend the control flow relationship and the intended sequence of information in code. Specifically, we randomly sample 20% of consecutive instruction pairs for example  $(I_1, I_2)$  and swap its sequence to be  $(I_2, I_1)$ , and then ask WasmRev to predict whether the instructions are reordered or not. We show an example of RII in Figure 7, where the instructions “local.tee 1” and “i32.load8\_u” are swapped. The loss function of RII objective is defined as:

$$\mathcal{L}_{RRI} = - \sum_{I_1, I_2 \in I} [y_{1,2} \log P(y_{1,2} | I_1, I_2) + (1 - y_{1,2}) \log(1 - P(y_{1,2} | I_1, I_2))] \quad (4)$$

where  $I$  is a set of instruction candidates for RRI,  $(I_1, I_2)$  is an instruction pair as input. The probability  $P(y_{1,2} | I_1, I_2)$  of instruction reordering is derived from the representation of the [CLS] token from WasmRev, following a softmax function.



Figure 7: An reordered instructions identification example

Finally the pre-training objective is to optimize WasmRev to minimize the loss function  $\mathcal{L}$  combining three tasks' loss functions:

$$\mathcal{L} = \mathcal{L}_{M3MLM} + \mathcal{L}_{SSI} + \mathcal{L}_{RII} + \lambda ||W||^2 \quad (5)$$

where  $W$  contains all trainable parameters of the model.  $\lambda$  is the coefficient of  $L_2$  regularizer to prevent overfitting.

## 5 DEVELOPING WEBASSEMBLY REVERSE ENGINEERING TOOLS

By exploiting the generalization of WasmRev and knowledge acquired from pre-training, we develop three WebAssembly reverse engineering tools for function purpose identification (FPI), type recovery (TR) and WebAssembly summarization (WS) by fine-tuning the pre-trained WasmRev correspondingly. These reverse engineering tasks take only the WebAssembly bytecode as input to generate high-level semantics, such as high-level types and function purpose, thereby facilitating the comprehension of WebAssembly programs.

### 5.1 Function Purpose Identification (FPI)

Function purpose identification is a task to identify the purpose of a given WebAssembly function, e.g., numeric processing, cryptography, etc. It provides developers with the type of a function, aiding in code understanding and further code inspection. We use the same problem setting as described in WASPur [61] and formulate FPI as a classification task to identify the type of a function from a set of predefined types. By leveraging the learnt WebAssembly representation in pre-training, we fine-tune WasmRev to accurately



predict function types. Specifically, we use the last hidden state of [CLS] token for classification to fine-tune WasmRev.

## 5.2 Type Recovery (TR)

Type recovery is a task to generate high-level precise types for parameters and return value of a function. It is an important first step toward WebAssembly comprehension, and is widely explored in existing reverse engineering tools for native binaries [9, 12, 58]. We employ the same problem setting as described in SnowWhite [42] and treat the type recovery problem as a sequence prediction task, where the sequence comprises tokens in a high-level type language. We leverage the knowledge of source code-WebAssembly relationship learnt in pre-training, and fine-tune WasmRev to recover parameters and return types. We employ an encoder-decoder model for this Wasm-to-PL generation task. Specifically, we use the pre-trained WasmRev to initialize the encoder and use a simple 2-layer bi-LSTM with an attention mechanism as the decoder.

## 5.3 WebAssembly Summarization (WS)

Code documentation is a crucial task to bridge the gap between complex code and understandable documentation, facilitating easier code comprehension, maintenance, and collaboration. Code documentation has been widely studied for source code [5, 21, 47]. Yet, the specific task of WebAssembly summarization has received limited attention. Existing tools like WASPur, which predicts the general purpose of functions, offering a rudimentary form of summarization that can be treated as a preliminary summarization. However, it falls short in capturing the specific functionalities of the code. For instance, WASPur might broadly classify a function as ‘numerical processing’ rather than precisely stating it ‘return the sum of two numbers.’, necessitating further code inspection for concrete understanding. To the best of our knowledge, this is the first work to formally define the WebAssembly summarization and propose a ML solution for it. We define WebAssembly summarization as a sequence generation task to generate human-readable natural language descriptions or summaries of a given WebAssembly code, i.e., what a particular piece of WebAssembly code does. This requires the model to understand not just the syntax but also the semantics of the WebAssembly code. A good summarization should be both factually correct, and useful and readable to the end-user. For WS, since it is a Wasm-to-NL generation task, we employ the same encoder-decoder architecture as utilized in the TR task, and fine-tune it on WS dataset.

# 6 IMPLEMENTATION

## 6.1 Pre-training Dataset

To facilitate generalizability of WasmRev and effective representation learning in pre-training, we collect a large scale of C/C++ code paired with documentation, from 432 publicly available open-source non-fork GitHub repositories, and prefer the high-quality projects indicated by their number of stars and forks. To obtain a representative dataset that has a close distribution to the realistic WebAssembly distribution, we follow the common WebAssembly use cases [16] to collect a broad spectrum of applications spanning diverse domains, e.g., gaming, text/media/numerical processing,

machine learning, cryptography, etc., and covering varying complexities and code styles.

We remove duplicates from the dataset by identifying (near) duplicate functions and only keeping one copy of them, following the de-duplication method proposed by Allamanis [4]. For the documentation, we truncate it to the first full paragraph, which typically summarizes and describes a function, while removing in-depth discussion of function arguments and return values. Samples with documentations shorter than three tokens are filtered out, as they are not informative. After preprocessing and filtering the source code and documentation, we collect 378k C/C++ functions paired with their documentation. We allocate only 80% of them as pre-training data of WasmRev, while the remaining 20% (75.6k) are used for the WS task, ensuring that WasmRev has never seen the source code and documentation in the WS task. We then compile the source code using Emscripten [15] v3.1.12 with different optimization options -O0, -O1, -O2, -O3, -Os, -Oz, allowing WasmRev to learn WebAssembly code with various compilations. We employ WABT [26] to transform WebAssembly into the text format. During this process, invalid samples that cannot be compiled by Emscripten or transformed by WABT are filtered out. For code file consisting of multiple functions, we first obtain the unstripped binary to identify the offsets of functions, and then use the offsets to locate the functions in the stripped binary and label them with the corresponding documentation. Ultimately, we collect ~ 1.8M (code documentation, source code, WebAssembly code) samples for WasmRev pre-training.

## 6.2 Fine-tuning Datasets

**6.2.1 FPI dataset:** We use the dataset from WASPur [61], consisting of 1,829 WebAssembly files collected from real-world websites, Firefox add-ons, Chrome extensions, and GitHub repositories, with 151,662 functions manually classified into 12 categories (listed in Figure 8).

**6.2.2 TR dataset:** We use the dataset from SnowWhite [42], composed of 5.5 million parameter type samples and 796 thousand return type samples, labeled with high-level types in its proposed expressive type language. It should be noted that the utilization of source code containing high-level types in C/C++ in pre-training does not result in data leakage for the TR task. This is because the TR task employs the type language defined in SnowWhite, making the targeted type sequences distinct from the C/C++ types used in pre-training data.

**6.2.3 WS dataset:** As mentioned in Section 6.1, we use 75.6k C/C++ code with documentation labels (which the pre-trained WasmRev has never seen) for WS task, and compile them with Emscripten and six various optimization options. We finally have 453k (WebAssembly code, documentation) samples for WS task.

All the datasets for fine-tuning tasks are split by 80%: 10%:10% into training / validation / testing sets. For data splitting, we split the source code first, then compile them using different optimization settings. In this way, we ensure functions from one project are not split over training, testing, and validation sets.

### 6.3 WasmRev Implementation

We implement our WasmRev model using PyTorch v1.11 [13] in python. For WasmRev’s neural architecture, we apply 8 Transformer layers, with 128 hidden states, 8 attention heads, initialized with Xavier [25]. We set the max sequence length as 512 and use the Adam optimizer [36] with a learning rate of 0.0005,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , weight decay rate of 0.01, and linear decay of the learning rate. To mitigate the over-fitting problem, we employ a dropout technique with a dropping ratio of 0.1. We pre-train WasmRev model with a batch size of 32 for 5 epochs using the pre-training dataset as described in Section 6.1. Subsequently, we fine-tune it with the same batch size for 2 epochs on each downstream task and set the learning rate to 3e-5, utilizing the fine-tuning dataset outlined in Section 6.2. During the fine-tuning phase for each downstream task, we fine-tune the pre-trained WasmRev model using the corresponding training set, and monitor accuracy on the validation set to facilitate early stopping upon model convergence. The best-performing model from the validation phase is then selected for the final evaluation using the testing set. All final predictions are obtained on the test data, which the model has never seen and was not used to select the best model. All the experiments are run on a Linux server running Ubuntu 20.04 with Intel Xeon E5-2686 v4 CPU with 12 cores running at 2.3Ghz, 488GB memory, and 8 Nvidia V100 16G GPUs.

## 7 EVALUATION

In this section, we report and analyze the experimental results to answer the following research questions (RQ):

- **RQ1:** Can WasmRev accurately distinguish different semantics of WebAssembly?
- **RQ2:** Can WasmRev accurately learn the relationship of WebAssembly and high-level programming language?
- **RQ3:** Can WasmRev accurately grasp the relationship of WebAssembly and natural language?
- **RQ4:** How effective are the multiple modalities of input representations?
- **RQ5:** How effective are our designed pre-training tasks?
- **RQ6:** How is the data-efficiency of WasmRev fine-tuning?

### 7.1 Function Purpose Identification (RQ1)

To study whether WasmRev can accurately distinguish between the various semantics of WebAssembly (RQ1), we fine-tune the pre-trained WasmRev model and evaluate it on the function purpose identification (FPI) task as described in Section 5.1, using the fine-tuning dataset specified in Section 6.2.1. FPI is a WebAssembly classification task, where the input is a WebAssembly function and output is the purpose of the input function, e.g., numeric processing, cryptography. The FPI task evaluates WasmRev’s ability on understanding and recognizing various semantic patterns of WebAssembly functions by requiring it to categorize these functions based on their specific purposes.

- **Baseline.** We use the SOTA model WASPur as our baseline for FPI, which employs a DNN model with an embedding layer, three hidden layers and an output layer, with supervised learning.
- **Metrics.** We adopt four metrics used in baseline WASPur: accuracy (Acc), precision (P), recall (R), F1 score. They are calculated by

true positive (TP), true negative (TN), false positive (FN) and false negative (FN), where true or false identify whether the prediction is correct or not, and positive and negative indicate the positive class or negative class. The metrics are defined as:

$$Acc = \frac{TP + TN}{TP + TN + FP + FN}, P = \frac{TP}{TP + FP}, R = \frac{TP}{TP + FN}, F1 = \frac{2 * P * R}{P + R} \quad (6)$$

- **Results.** We show the results in Table 1. WasmRev outperforms WASPur by 6.18%, 3.30%, 5.98%, 4.61% in terms of top-1 accuracy, precision, recall and F1 score, revealing the effectiveness of WasmRev to learn WebAssembly representations, and recognize diverse semantic patterns among various types of functions.

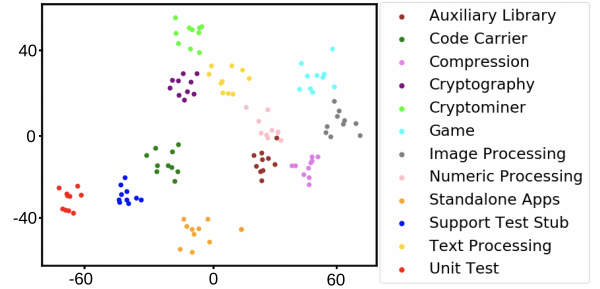
**Table 1: FPI and WS results, and effectiveness of design components. WasmRev and its variants are fine-tuned separately for various tasks.**

Task	FPI				TR	TR	WS		
	Metric	Acc*	P	R	F1	param	return	BLEU	BF1 <sup>‡</sup>
WASPur <sup>†</sup>	88.07%	0.91	0.87	0.89	-	-	-	-	-
WasmRev <sub>SuV(1/10)</sub> (WS)	-	-	-	-	-	-	-	18.62	0.765
WasmRev <sub>SuV</sub> (WS)	-	-	-	-	-	-	-	19.57	0.833
WasmRev <sub>w/oNL</sub>	93.17%	0.935	0.919	0.927	63.2%	72.5%	19.88	0.840	
WasmRev <sub>w/oPL</sub>	93.20%	0.933	0.917	0.925	62.8%	72.1%	19.57	0.831	
WasmRev <sub>w/oNL+PL</sub>	92.88%	0.921	0.910	0.915	62.5%	71.9%	19.49	0.828	
WasmRev <sub>w/oM3MLM</sub>	93.06%	0.930	0.915	0.922	62.2%	71.8%	20.00	0.851	
WasmRev <sub>w/oSSI</sub>	93.00%	0.929	0.916	0.922	62.8%	72.2%	20.10	0.855	
WasmRev <sub>w/oRII</sub>	92.92%	0.925	0.908	0.916	62.6%	72.1%	20.12	0.855	
WasmRev <sub>(1/10)</sub>	91.25%	0.920	0.901	0.910	61.9%	70.5%	19.32	0.831	
WasmRev	<b>93.51%</b>	<b>0.940</b>	<b>0.922</b>	<b>0.931</b>	<b>63.7%</b>	<b>72.9%</b>	<b>20.33</b>	<b>0.873</b>	

\* “Acc” indicates top-1 accuracy.

<sup>†</sup> The best results reported in the WASPur paper.

<sup>‡</sup> “BF1” indicates BERTScore F1 score.



**Figure 8: Visualization of FPI results of WasmRev.**

- **Visualization.** To intuitively show how well does WasmRev distinguish different semantics of WebAssembly (RQ1), we visualize the classification results. In particular, we randomly select 10 functions from each class of function purpose, and compute the embeddings of them. We then use the t-distributed Stochastic Neighbor Embedding (t-SNE) method [67] to project the high-dimensional embeddings onto a 2-D plane. We show the visualization in Figure 8, with different classes marked in different colors. The clear clustering of functions of the same type and the separation between functions of different types, reveal WasmRev’s capability of identifying and distinguishing various types of functions. Interestingly, we observe the representations of *Standalone Apps* are relative scattering, since various Apps could have diverse functionalities. One potential improvement is to refine the dataset, e.g., divide samples



in *Standalone Apps* into sub-classes of similar functionalities. We also observe some clusters are close to each other, e.g., *Numeric Processing* and *Auxiliary Library*, which could both involve utility functions, possibly with a mathematical focus.

Init_circles.wat	
...	-Source code piece:
f64.promote_f32	circleData[i].vy=(randomf() - 0.5)*0.01;
f64.const -0x1p-1 ;;=-0.5	
f64.add	
f64.const 0x1.47ae147ae147bp-7 ;;=0.03	-WasmRev's Prediction:
f64.mul	Numerical Processing
f32.demote_f64	
f32.store	-GroundTruth:
...	Auxiliary Library

Figure 9: A case study of FPI.

• **Case Studies.** We show two case studies of FPI. Since we do not have access to the source code of WASPur, we limit our discussion to our results. The first case is shown in Figure 3, where WasmRev accurately identifies the function type as text processing. The second case, illustrated in Figure 9, is a function to initialize the positions and velocities of circles to be used in a graphical application, belong to the *Auxiliary Library* type. Here, WasmRev incorrectly categorizes the function type as *Numerical Processing*, since the function involves a lot of mathematical operations as shown in the WebAssembly code `Init_circles.wat`. This aligns with the observation in the visualization that the ambiguity of code and labeling could lead to incorrect predictions. To address this issue, we could use the results from the WS task to aid in better comprehending the function purpose.

**Result-1:** WasmRev is able to accurately distinguish different semantics of WebAssembly, with a high accuracy (93.51%) to identify the correct class of purpose of a WebAssembly function, outperforming WASPur by 6.18% in top-1 classification accuracy.

## 7.2 Type Recovery (RQ2)

To investigate whether WasmRev can accurately learn the relationship of WebAssembly and high-level programming language (RQ2), we fine-tune the pre-trained WasmRev model and evaluate it on the type recovery (TR) task as detailed in Section 5.2, with the fine-tuning dataset described in Section 6.2.2. TR is a sequence generation task, where the input is a WebAssembly function and output is a sequence in high level programming language that specifies the parameters and return types, e.g., `primitive int32`, `pointer primitive char`. TR evaluates how well WasmRev can map concepts and types in low-level WebAssembly to high-level type signatures by understanding the WebAssembly context.

• **Baseline.** We use the SOTA SnowWhite as our baseline, which uses supervised learning on a sequence-to-sequence model [64].

• **Metrics.** We adopt three metrics used in baseline SnowWhite. (1) top-1 accuracy (perfect match accuracy). (2) top-5 accuracy, which retrieves the five most likely type predictions via beam search. (3) The average type prefix score (TPS) over the whole test set, where the TPS is defined as the length of the common prefix  $TPS(t', t) = |commonPrefix(t', t)|$  of a prediction  $t'$  and ground truth  $t$ .

• **Results.** We show the results in Table 2, with various high-level type languages provided in the SnowWhite dataset. With the  $\mathcal{L}_{SW}$  type language and 1,225 unique types in the testing set, WasmRev achieves 63.7%/88.3% top-1/top-5 accuracy for parameter type

Top-5 predictions for the return type of function:  
`int JPEGVGetField(TIFF* tif, uint32 tag, va_list ap) ...`

SnowWhite	WasmRev
primitive int 8	<b>primitive int 32</b>
primitive uint 32	primitive uint 32
primitive uint 32	pointer primitive int 32
pointer name "Exception" class	primitive int 8
<b>primitive int 32</b>	pointer primitive uint 32

Figure 10: A case study of TR.

prediction and 74.9%/93.8% top-1/top-5 accuracy for return type prediction, outperforming SnowWhite by 43.1%/17.4% and 29.8%/16.5% correspondingly. WasmRev achieves an average type prefix score of 1.78 for parameters and 1.60 for return values. This indicates that WasmRev can accurately predict the first 1.78 and 1.60 tokens of the type sequence for parameters and return values, respectively. The results reveal WasmRev's effectiveness in type recovery, since we expect the one or two tokens to be likely the most relevant to a reverse engineer. ( $\mathcal{L}_{SW}$ , All Names) is a variant of  $\mathcal{L}_{SW}$  with non-filtered type names and 146,883 different types in total, making the TR task more difficult. ( $\mathcal{L}_{SW}$ , Simplified) is a simplified variant of  $\mathcal{L}_{SW}$ , with some constructors removed in the type language to consist of 120 types.  $\mathcal{L}_{Eklavya}$  has the simplest setting with only 7 different types. ( $\mathcal{L}_{SW}$ ,  $t_{low}$  not given) is similar to  $\mathcal{L}_{SW}$ , but removes the low-level type information in the input. Among all these settings, WasmRev achieves better results than SnowWhite across all the metrics, demonstrating its effectiveness in understanding the relationship between WebAssembly and high-level programming language, and projecting the correct types from WebAssembly.

• **Case Studies.** We show two case studies of TR. The first involves simple string manipulation, as shown in Figure 3. Here WasmRev accurately identifies the parameter type as `pointer primitive char`. The second case, illustrated in Figure 10, showcases the top five predictions for the return type of the `JPEGVGetField` function within the `libtiff` library. In the results from SnowWhite, the correct prediction ranks fifth. In contrast, WasmRev precisely predicts the return type as `primitive int 32` in its top-1 prediction, and the other predictions within the top-5 are also close to the ground truth.

**Result-2:** WasmRev is able to accurately learn the relationship of WebAssembly and high-level programming language, and forecast precise parameter and return types from WebAssembly with high top-1/top-5 accuracy 63.7%/88.3% and 74.9%/93.8%, outperforming SnowWhite by 43.1%/17.4% and 29.8%/16.5% correspondingly.

## 7.3 WebAssembly Summarization (RQ3)

To explore whether WasmRev can precisely understand the relationship of WebAssembly and natural language (RQ3), we fine-tune the pre-trained WasmRev model and evaluate it on the WebAssembly summarization (WS) task as introduced in Section 5.3, using the fine-tuning dataset detailed in Section 6.2.3. WS is a sequence generation task, where the input is a WebAssembly function and output is a sequence of natural language summaries describing that function, thereby testing model's capacity to bridge the gap between complex, low-level code constructs and their high-level, linguistic descriptions.

• **Baseline.** As there is no existing learning-based solution for WebAssembly summarization, we train a supervised learning model

**Table 2: TR results.**

Task		Parameter Type Prediction					Return Type Prediction				
Type	Language	$\mathcal{L}_{SW}$	$\mathcal{L}_{SW}$ , All Names	$\mathcal{L}_{SW}$ , Simplified	$\mathcal{L}_{Eklavya}$	$\mathcal{L}_{SW}$ , $t_{low}$ not given	$\mathcal{L}_{SW}$	$\mathcal{L}_{SW}$ , All Names	$\mathcal{L}_{SW}$ , Simplified	$\mathcal{L}_{Eklavya}$	$\mathcal{L}_{SW}$ , $t_{low}$ not given
SnowWhite <sup>†</sup>	Top-1 Acc	44.5%	18.6%	65.1%	87.9%	43.4%	57.7%	40.6%	60.6%	76.3%	50.7%
	Top-5 Acc	75.2%	27.1%	86.2%	100.0%	74.3%	80.5%	47.3%	87.9%	100%	81.2%
	Type Prefix Score	1.47	1.31	1.62	0.88	1.45	1.37	1.00	1.38	0.76	1.02
WasmRev	Top-1 Acc	<b>63.7%</b>	40.2%	80.6%	93.4%	62.8%	<b>74.9%</b>	52.4%	79.7%	89.2%	73.5%
	Top-5 Acc	88.3%	50.0%	95.2%	100.0%	87.4%	93.8%	63.9%	95.1%	100%	93.2%
	Type Prefix Score	1.78	1.55	1.89	0.93	1.76	1.60	1.30	1.58	0.89	1.31

<sup>†</sup> The results reported in the SnowWhite paper.

with the same architecture as WasmRev on the WS dataset as our baseline, referred to as WasmRev<sub>SuV</sub>(WS).

• **Metrics.** We use the widely-used bilingual evaluation understudy (BLEU) score [45], and BERTScore F1 score [75] to compare the generated sentence to reference sentence. Different from the normal F1 score, BERTScore utilizes the contextual embeddings from the BERT model, and compares the semantic similarity of texts by calculating the cosine similarity between the embeddings of the generated and reference texts. The higher BLEU score/BERTScore F1 score corresponds to a generated sentence that is more similar to the reference.

• **Results.** We show the results in Table 1. WasmRev outperforms the supervised learning baseline by 3.9% and 4.8% in BLEU-4 score and BERTScore F1 score, revealing WasmRev’s superior ability on understanding WebAssembly and describing it in human-readable summaries.

• **Case Studies.** We present two case studies of WS. One is a simple string manipulation shown in Figure 3, where the WS result of WasmRev is most the same as the ground-truth documentation and presents the correct semantics. In this example, WasmRev<sub>SuV</sub>(WS) generates almost the same summarization: “Convert the lowercase letters to uppercase”, which is also semantically correct. Another case is shown in Figure 11 with partial WebAssembly code. It is a more complex hash function combining bit-wise operations and arithmetic operations, with a very detailed documentation describing its functionality. WasmRev precisely predicts the function purpose – “a hash function”, and correctly includes operations like “bit shifts” and “additions”. It incorrectly includes “multiplications”. We guess a possible reason is the appearance of the “i32.mul” between parameters 1 and 2 (“local.get 1/2”) in the WebAssembly code (as shown in Figure 11) makes the model predict the “multiplications”. The supervised learning baseline WasmRev<sub>SuV</sub>(WS) also forecasts the “shift” and “accumulate” operations but with the wrong object – it is the hash value that is bit shifted, not the characters of the string. In this case, WasmRev predicts more meaningful tokens and produces a semantically closer summarization to the ground-truth.

<pre>;;... loop   local.get 2   i32.const 65599   i32.mul   local.get 1   i32.const 24   i32.shl   i32.const 24   i32.shr_s   i32.add ;;...</pre>	<p><b>GroundTruth:</b> Calculate a hash value of a string by iteratively combining each character with the current hash value using bit shifts and additions.</p> <p><b>WasmRev:</b> A hash function with bit shifts, multiplications and additions to accumulate the hash value.</p> <p><b>WasmRev_SuV(WS):</b> Shift each character of a given string and accumulate the sum.</p>
---	---

**Figure 11: A case study of WS.**

**Result-3:** WasmRev is able to accurately summarize the semantics of WebAssembly in natural language, with a high BLEU-4 score (20.33) and a high BERTScore F1 (0.873) compared to non-pretrained model.

## 7.4 Effectiveness of multi-model inputs (RQ4)

To evaluate the effectiveness of using multi-modal inputs, we train the following variants for each task: (1) WasmRev<sub>w/oNL</sub> uses only (PL, Wasm) as input and accordingly removes the NL samples in SSI task; (2) WasmRev<sub>w/oPL</sub> uses only (NL, Wasm) as input and accordingly removes the PL samples in SSI task; (3) WasmRev<sub>w/oNL+PL</sub> uses only WebAssembly as input and uses only the different compiled versions of WebAssembly code as the similar semantics in SSI task. Other settings are identical to WasmRev. The results are shown in Table 1 for FPI, TR and WS tasks.

WasmRev outperforms all the variants for three tasks, indicating the effectiveness of each modality of input representation. We observe removing various modalities affect differently for different tasks. (i) For FPI, removing NL, PL both leads to performance degradation while removing NL+PL incurs the biggest accuracy drop (-0.63%), indicating incorporating NL and PL are both important and helpful for WasmRev to learn WebAssembly semantics. (ii) For TR, removing PL leads to a greater performance degradation (-0.9%/-0.8%) than removing only NL (-0.5%/-0.4%), for parameter/return TR tasks, indicating learning PL and WebAssembly concurrently helps modeling PL-WebAssembly relationship, contributing to a promising result of Wasm-to-PL task. (iii) For WS, removing NL leads to a greater performance degradation (-0.76) than removing only PL (-0.45), revealing that co-learning NL and Wasm helps the model grasp high-level semantics and NL-Wasm relationship.

**Result-4:** Our results empirically demonstrate that incorporating various modalities of code into WasmRev contributes to improved model performance for various reverse engineering tasks.

## 7.5 Effectiveness of pre-training tasks (RQ5)

Pre-training is a critical process in which WasmRev learns multi-modal representations and cross-modal relationships. To evaluate the effectiveness of our designed pre-training tasks, we train three kinds of variants, each of which removes one of the pre-training task T, denoted as WasmRev<sub>w/oT</sub>, T ∈ {M3MLM, SSI, RII}.

The results are shown Table 1, where WasmRev outperforms all the variants, showing the effectiveness of each of our design pre-training tasks. Specifically, (i) For FPI, removing any of the pre-training task incur performance degradation, while removing

RII has a relative greater affect to FPI task. (ii) For TR, removing M3MLM leads to the biggest accuracy drop. The M3MLM task, which involves predicting masked source code based on its context and the semantically equivalent WebAssembly, effectively enables the model to learn the PL-Wasm relationship. (iii) For WS, similarly, M3MLM brings the greatest impact since the model learns crucial NL-Wasm relationship through M3MLM, which is valuable for WS.

**Result-5:** We empirically show that our designed pre-training tasks yield promising results for various WebAssembly tasks.

## 7.6 Data Efficiency (RQ6)

To evaluate the data efficiency of fine-tuning, i.e., the ability to achieve high model performance with a smaller amount of labeled data, we train WasmRev variants using only  $\frac{1}{10}$  of labeled data for WasmRev fine-tuning. We show the results in Table 1. For FPI and TR, although WasmRev<sub>(1/10)</sub> exhibits a minor reduction in performance compared to WasmRev, it still surpasses baselines that utilize the full dataset for training. For WS, WasmRev<sub>(1/10)</sub> attains a BLEU-4 score and a BERTScore F1 that are both nearly on par with, yet marginally below that of WasmRev<sub>SuV</sub>(WS). Notably, WasmRev<sub>(1/10)</sub> yields a higher BLEU-4 and a higher BERTScore F1 than WasmRev<sub>SuV</sub>(1/10)(WS) that uses the same  $\frac{1}{10}$  of labeled data for supervised training. These findings underscore WasmRev's capacity to perform efficiently in scenarios where obtaining labeled data is either challenging or resource-intensive, thereby addressing the limitation L1 outlined in Section 2.3 and highlighting the practical advantages of WasmRev in data-constrained environments.

**Result-6:** Even with significantly less (1/10) labeled data for fine-tuning, WasmRev still achieves high model performance and outperforms the SOTA WASPur and SnowWhite.

## 7.7 Time Efficiency

We discuss the time efficiency across three stages of WasmRev. (1) The model pre-training takes  $\sim 1.5$  days on GPUs. Note that the training efficiency can be further improved by using more advanced machines such as Nvidia A100 or H100 GPUs [53, 55], and training techniques such as mixed precision training [56], with a reported 1.2 hours for BERT model training [54]. The time cost of pre-training is a *one-time* overhead, which is more efficient than conventional WebAssembly analysis tool development that can often take expert engineers a few weeks or months. It is more efficient than supervised training, which requires more labeled data collection effort and fully training of each model for individual tasks, resulting in a total time overhead scaling linearly with the number of tasks. In this sense, WasmRev improves WebAssembly tool development efficiency compared to conventional methods and supervised learning methods. (2) In the fine-tuning stage, it takes 0.5 / 9 / 1.5 hours to fine-tune WasmRev on FPI / TR / WS task. This overhead depends greatly on the size of fine-tuning dataset. (3) The model inference takes on average 350ms per input binary in our testing, depending on the input sequence length. Such fast inference further enhances the practicality of these reverse engineering tools developed based on WasmRev.

## 8 RELATED WORK

**Conventional WebAssembly analysis.** Since the introduction of the WebAssembly standard [1, 29], several studies have been conducted for WebAssembly security [37, 40, 51], performance [33], and use in practice [30]. Techniques to analyze and inspect WebAssembly code include static analysis [8], dynamic analysis [41] and taint analyses [22, 65]. WABT [26] supports converting WebAssembly code to C source code and decompiling WebAssembly code into C-like syntax. However, none of them targets high-level semantics recovery, e.g., types, function purpose and code documentation.

**ML for WebAssembly reverse engineering.** To the best of our knowledge, there are very few previous works [42, 61] targeting ML for WebAssembly reverse engineering. SnowWhite [42] learns WebAssembly sequences and focuses on type recovery. WASPur [61] extracts control-flow graphs of WebAssembly and targets function purpose classification. These ML methods focus on specific applications while WasmRev learns a generic WebAssembly representation for reverse engineering tasks. Recent SOTA large language models (LLMs) such as ChatGPT [57] are not designed for WebAssembly analysis and commercial LLMs are not always viable options for WebAssembly analysis due to potential privacy concerns. To the best of our knowledge, WasmRev is the first pre-trained language model for generalized WebAssembly reverse engineering tasks.

**Neural models of code.** DL techniques have been extensively explored for code-related tasks, such as code summarization [39, 68], vulnerability detection [22, 44, 65] and binary similarity comparison [46, 66, 74]. One avenue is to learn code representation, i.e., learn to estimate the statistical distributional properties over large and representative corpora. Earlier approaches apply word2vec [50] to obtain code representations [7, 11]. More recent works exploit representation learning on SQL query [32], source-level programming languages [17, 27], binary code [58, 69], and bi-modal (NL-PL) [21, 28]. Our exploration resonates with recent studies on the principle of neural software analysis [60] and code representation learning [52], while proposing the first multi-modal representation learning of NL-PL-Wasm for generalized WebAssembly tasks.

## 9 CONCLUSION

This paper proposes WasmRev, the first multi-modal pre-trained language model, which learns a generic representation of source code, code documentation, and WebAssembly code. It can be applied to diverse WebAssembly reverse engineering tasks. Incorporating multi-modal inputs and tailored pre-training tasks, WasmRev can learn from cross-modal relationships and gain a holistic understanding of WebAssembly. WasmRev generates accurate results for type recovery, function purpose identification, and WebAssembly summarization. It outperforms supervised learning methods, providing developers with high-level semantics for inspecting WebAssembly modules and assisting WebAssembly comprehension. We leave exploring WasmRev for other reverse engineering tasks such as function name prediction and precise decompilation as future work.

## 10 ACKNOWLEDGMENT

The authors thank the anonymous reviewers for their valuable feedback and comments. This paper is supported in part by NSF grants 1829524, 1817077, 2011212, and the PRISM center in JUMP 2.0, an SRC program sponsored by DARPA.



## REFERENCES

- [1] 2018. WebAssembly website. <https://webassembly.org/>
- [2] 2020. GitHub. <https://github.com>
- [3] 2022. Wasm non web usage. <https://webassembly.org/docs/non-web/>
- [4] Miltiadis Allamanis. 2019. The adverse effects of code duplication in machine learning models of code. In *Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*. 143–153.
- [5] Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav. 2018. code2seq: Generating sequences from structured representations of code. *arXiv preprint arXiv:1808.01400* (2018).
- [6] Nuttapong Attrapadung, Goichiro Hanaoka, Shigeo Mitsunari, Yusuke Sakai, Kana Shimizu, and Tadanori Teruya. 2018. Efficient two-level homomorphic encryption in prime-order bilinear groups and a fast implementation in web-assembly. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*. 685–697.
- [7] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the association for computational linguistics* 5 (2017), 135–146.
- [8] Tiago Brito, Pedro Lopes, Nuno Santos, and José Frago Santos. 2022. Wasmati: An efficient static vulnerability scanner for WebAssembly. *Computers & Security* 118 (2022), 102745.
- [9] Juan Caballero and Zhiqiang Lin. 2016. Type inference on executables. *ACM Computing Surveys (CSUR)* 48, 4 (2016), 1–35.
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [11] Zimin Chen and Martin Monperrus. 2019. A literature study of embeddings on source code. *arXiv preprint arXiv:1904.03061* (2019).
- [12] Zheng Leong Chua, Shiqi Shen, Prateek Saxena, and Zhenkai Liang. 2017. Neural nets can learn function type signatures from binaries. In *26th USENIX Security Symposium (USENIX Security 17)*. 99–116.
- [13] Roman Collobert, Samy Bengio, and Johnny Mariéthoz. 2002. *Torch: a modular machine learning software library*. Technical Report. Idiap.
- [14] DWARF Committee. 2017. DWARF Debugging Information Format – Version 5. <http://www.dwarfstd.org/doc/DWARF5.pdf>
- [15] Emscripten Contributors. 2020. Emscripten. <https://emscripten.org/index.html>
- [16] WebAssembly Contributors. 2020. Webassembly Use Cases. <https://webassembly.org/docs/use-cases/>
- [17] Hoa Khanh Dam, Truyen Tran, and Trang Pham. 2016. A deep language model for software code. *arXiv preprint arXiv:1608.02715* (2016).
- [18] Nikhil Thorat Daniel Smilov and Ann Yuan. 2020. Introducing the WebAssembly backend for TensorFlow.js. <https://blog.tensorflow.org/2020/03/introducing-webassembly-backend-for-tensorflow-js.html>
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [20] Yulong Wang Emma Ning and Du Li. 2021. ONNX Runtime Web—running your machine learning model in browser. <https://cloudblogs.microsoft.com/opensource/2021/09/02/onnx-runtime-web-running-your-machine-learning-model-in-browser/>
- [21] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).
- [22] William Fu, Raymond Lin, and Daniel Inge. 2018. Taintassembly: Taint-based information flow control tracking for webassembly. *arXiv preprint arXiv:1802.01050* (2018).
- [23] Phani Kishore Gadepalli, Sean McBride, Gregor Peach, Ludmila Cherkasova, and Gabriel Parmer. 2020. Sledge: A serverless-first, light-weight wasm runtime for the edge. In *Proceedings of the 21st International Middleware Conference*. 265–279.
- [24] Phani Kishore Gadepalli, Gregor Peach, Ludmila Cherkasova, Rob Aitken, and Gabriel Parmer. 2019. Challenges and opportunities for efficient serverless computing at the edge. In *2019 38th Symposium on Reliable Distributed Systems (SRDS)*. IEEE, 261–2615.
- [25] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 249–256.
- [26] W. C. Group. 2019. webassembly/wabt. <https://github.com/WebAssembly/wabt>
- [27] Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. 2016. Deep API learning. In *Proceedings of the 2016 24th ACM SIGSOFT international symposium on foundations of software engineering*. 631–642.
- [28] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366* (2020).
- [29] Andreas Haas, Andreas Rossberg, Derek L Schuff, Ben L Titzer, Michael Holman, Dan Gohman, Luke Wagner, Alon Zakai, and JF Bastien. 2017. Bringing the web up to speed with WebAssembly. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*. 185–200.
- [30] Aaron Hilbig, Daniel Lehmann, and Michael Pradel. 2021. An empirical study of real-world webassembly binaries: Security, languages, use cases. In *Proceedings of the web conference 2021*. 2696–2708.
- [31] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. CodeSearchNet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436* (2019).
- [32] Shrainer Jain, Bill Howe, Jiaqi Yan, and Thierry Cruanes. 2018. Query2vec: An evaluation of NLP techniques for generalized workload analytics. *arXiv preprint arXiv:1801.05613* (2018).
- [33] Abhinav Jangda, Bobby Powers, Emery D Berger, and Arjun Guha. 2019. Not so fast: Analyzing the performance of {WebAssembly} vs. native code. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. 107–120.
- [34] Hunseop Jeong, Jinwoo Jeong, Sangyong Park, and Kwanghyuk Kim. 2018. WATT: A novel web-based toolkit to generate WebAssembly-based libraries and applications. In *2018 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 1–2.
- [35] Suvarna Kadam and Vinay Vaidya. 2020. Review and analysis of zero, one and few shot learning approaches. In *Intelligent Systems Design and Applications: 18th International Conference on Intelligent Systems Design and Applications (ISDA 2018) held in Vellore, India, December 6–8, 2018, Volume 1*. Springer, 100–112.
- [36] Diederik Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*. San Diego, CA, USA.
- [37] Radhesh Krishnan Konoth, Emanuele Vineti, Veelasha Moonsamy, Martina Lindorfer, Christopher Kruegel, Herbert Bos, and Giovanni Vigna. 2018. Minesweeper: An in-depth look into drive-by cryptocurrency mining and its defense. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 1714–1730.
- [38] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942* (2019).
- [39] Alexander LeClair, Siyuan Jiang, and Collin McMillan. 2019. A neural model for generating natural language summaries of program subroutines. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 795–806.
- [40] Daniel Lehmann, Johannes Kinder, and Michael Pradel. 2020. Everything old is new again: Binary security of {WebAssembly}. In *29th USENIX Security Symposium (USENIX Security 20)*. 217–234.
- [41] Daniel Lehmann and Michael Pradel. 2019. Wasabi: A framework for dynamically analyzing webassembly. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 1045–1058.
- [42] Daniel Lehmann and Michael Pradel. 2022. Finding the dwarf: recovering precise types from WebAssembly binaries. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation*. 410–425.
- [43] Xuezixiang Li, Yu Qu, and Heng Yin. 2021. Palmtree: Learning an assembly language model for instruction embedding. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 3236–3251.
- [44] Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Yawei Zhu, and Zhaoxuan Chen. 2021. Sysevr: A framework for using deep learning to detect software vulnerabilities. *IEEE Transactions on Dependable and Secure Computing* 19, 4 (2021), 2244–2258.
- [45] Chin-Yew Lin and Franz Josef Och. 2004. Orange: a method for evaluating automatic evaluation metrics for machine translation. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*. 501–507.
- [46] Bingchang Liu, Wei Huo, Chao Zhang, Wencho Li, Feng Li, Aihua Piao, and Wei Zou. 2018.  $\alpha$ Diff: Cross-Version Binary Code Similarity Detection with DNN. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (Montpellier, France) (ASE '18)*. Association for Computing Machinery, New York, NY, USA, 667–678. <https://doi.org/10.1145/3238147.3238199>
- [47] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664* (2021).
- [48] Niko Mäkitalo, Tommi Mikkonen, Cesare Pautasso, Victor Bankowski, Paulius Daubaris, Risto Mikkola, and Oleg Beletski. 2021. WebAssembly modules as lightweight containers for liquid IoT applications. In *International Conference on Web Engineering*. Springer, 328–336.
- [49] Vadim Markovtsev and Warren Long. 2018. Public git archive: a big code dataset for all. In *Proceedings of the 15th International Conference on Mining Software Repositories*. 34–37.
- [50] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).

- [51] Marius Musch, Christian Wressnegger, Martin Johns, and Konrad Rieck. 2019. New Kid on the Web: A Study on the Prevalence of WebAssembly in the Wild. In *Detection of Intrusions and Malware, and Vulnerability Assessment: 16th International Conference, DIMVA 2019, Gothenburg, Sweden, June 19–20, 2019, Proceedings* 16. Springer, 23–42.
- [52] Changan Niu, Chuanyi Li, Bin Luo, and Vincent Ng. 2022. Deep learning meets software engineering: A survey on pre-trained models of source code. *arXiv preprint arXiv:2205.11739* (2022).
- [53] Nvidia. 2023. Accelerating the Most Important Work of Our Time. <https://www.nvidia.com/en-us/data-center/a100/>
- [54] Nvidia. 2023. DeepLearningExamples. <https://github.com/NVIDIA/DeepLearningExamples/blob/master/PyTorch/LanguageModeling/BERT/README.md#inference-performance-benchmark>
- [55] Nvidia. 2023. An Order-of-Magnitude Leap for Accelerated Computing. <https://www.nvidia.com/en-us/data-center/h100/>
- [56] Nvidia. 2023. Train With Mixed Precision. <https://docs.nvidia.com/deeplearning/performance/mixed-precision-training/index.html>
- [57] OpenAI. 2023. ChatGPT. <https://chat.openai.com/chat>
- [58] Kexin Pei, Jonas Guan, Matthew Broughton, Zhongtian Chen, Songchen Yao, David Williams-King, Vikas Ummadisetty, Junfeng Yang, Baishakhi Ray, and Suman Jana. 2021. 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*. 690–702.
- [59] Farhad Pourpanah, Moloud Abdar, Yuxuan Luo, Xinlei Zhou, Ran Wang, Chee Peng Lim, Xi-Zhao Wang, and QM Jonathan Wu. 2022. A review of generalized zero-shot learning methods. *IEEE transactions on pattern analysis and machine intelligence* (2022).
- [60] Michael Pradel and Satish Chandra. 2021. Neural Software Analysis. *Commun. ACM* 65, 1 (dec 2021), 86–96. <https://doi.org/10.1145/3460348>
- [61] Alan Romano and Weihang Wang. 2023. Automated WebAssembly Function Purpose Identification With Semantics-Aware Analysis. In *Proceedings of the ACM Web Conference 2023*. 2885–2894.
- [62] Alan Romano, Yunhui Zheng, and Weihang Wang. 2020. Minerray: Semantics-aware analysis for ever-evolving cryptojacking detection. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*. 1129–1140.
- [63] T. Rust and W. W. Group. 2020. wasm-bindgen. <https://github.com/rustwasm/wasm-bindgen>
- [64] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems* 27 (2014).
- [65] Aron Szanto, Timothy Tamm, and Artidoro Pagnoni. 2018. Taint tracking for webassembly. *arXiv preprint arXiv:1807.08349* (2018).
- [66] Donghai Tian, Xiaoqi Jia, Rui Ma, Shuke Liu, Wenjing Liu, and Changzhen Hu. 2021. BinDeep: A deep learning approach to binary code similarity detection. *Expert Systems with Applications* 168 (2021), 114348.
- [67] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [68] Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S Yu. 2018. Improving automatic source code summarization via deep reinforcement learning. In *Proceedings of the 33rd ACM/IEEE international conference on automated software engineering*. 397–407.
- [69] Hao Wang, Wenjie Qu, Gilad Katz, Wenyu Zhu, Zeyu Gao, Han Qiu, Jianwei Zhuge, and Chao Zhang. 2022. jtrans: Jump-aware transformer for binary code similarity. *arXiv preprint arXiv:2205.12713* (2022).
- [70] Wenhao Wang, Benjamin Ferrell, Xiaoyang Xu, Kevin W Hamlen, and Shuang Hao. 2018. Seismic: Secure in-lined script monitors for interrupting cryptojacks. In *Computer Security: 23rd European Symposium on Research in Computer Security, ESORICS 2018, Barcelona, Spain, September 3–7, 2018, Proceedings, Part II* 23. Springer, 122–142.
- [71] Conrad Watt, John Renner, Natalie Popescu, Sunjay Cauligi, and Deian Stefan. 2019. CT-Wasm: Type-Driven Secure Cryptography for the Web Ecosystem. *Proc. ACM Program. Lang.* 3, POPL, Article 77 (jan 2019), 29 pages. <https://doi.org/10.1145/3290390>
- [72] Elliott Wen and Gerald Weber. 2020. Wasmachine: Bring iot up to speed with a webassembly os. In *2020 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*. IEEE, 1–4.
- [73] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [74] Shouguo Yang, Long Cheng, Yicheng Zeng, Zhe Lang, Hongsong Zhu, and Zhiqiang Shi. 2021. Asteria: Deep learning-based AST-encoding for cross-platform binary code similarity detection. In *2021 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 224–236.
- [75] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675* (2019).

Received 16-DEC-2023; accepted 2024-03-02