

BinBert: Binary Code Understanding with a Fine-tunable and Execution-aware Transformer

Fiorella Artuso¹, Marco Mormando¹, Giuseppe Antonio Di Luna¹, and Leonardo Querzoni¹

¹Sapienza University of Rome

August 16, 2022

Abstract

A recent trend in binary code analysis promotes the use of neural solutions based on instruction embedding models. An instruction embedding model is a neural network that transforms sequences of assembly instructions into embedding vectors. If the embedding network is trained such that the translation from code to vectors partially preserves the semantic, the network effectively represents an *assembly code model*.

In this paper we present BinBert, a novel assembly code model. BinBert is built on a transformer pre-trained on a huge dataset of both assembly instruction sequences and symbolic execution information. BinBert can be applied to assembly instructions sequences and it is *fine-tunable*, i.e. it can be re-trained as part of a neural architecture on task-specific data. Through fine-tuning, BinBert learns how to apply the general knowledge acquired with pre-training to the specific task.

We evaluated BinBert on a multi-task benchmark that we specifically designed to test the understanding of assembly code. The benchmark is composed of several tasks, some taken from the literature, and a few novel tasks that we designed, with a mix of intrinsic and downstream tasks.

Our results show that BinBert outperforms state-of-the-art models for binary instruction embedding, raising the bar for binary code understanding.

1 Introduction

A recent body of works in the literature has shown that Deep Neural Networks (DNNs) can successfully solve several binary analysis tasks. DNNs today show state of the art performances for binary similarity [22, 42, 41], compiler provenance [8, 23, 32], function boundaries detection [33], decompiling [16], automatic function naming [11, 29] and others.

DNN designers must decide how to feed binary code to their models. One possibility is to use manually-identified features. This approach requires a domain expert which identifies features of interest forecasting their helpfulness in solving the task at hand. This approach is known to produce problem-specific features and injects a human bias inside the system. Recent solutions automatically transform binary code into a representation usable by the neural network layers.

A common technique is to transform assembly instructions into representational *embeddings vectors*, similarly to what has been done in the Natural Language Processing (NLP) field with the word embedding revolution [24]. Several works [22, 42, 9, 19] proposed refined techniques to transform a single instruction into a vector of real numbers while capturing its semantic (e.g. all vectors of arithmetic instructions are clustered in the vector

space). By using this approach, sequences of instructions are transformed into sequences of fixed-size vectors that can be fed into standard DNNs.

1.1 Execution-aware Binary Code Interpretation

A common weakness of all these solutions is *the lack of context*: an instruction is always represented by the same vector, irrespectively of where it appears. However, the semantic of a single assembly instruction is strongly limited (more than a word in natural language), and non-trivial concepts in assembly code are almost always encoded by a sequence of instructions (e.g., loops, swap of variables in memory, calling conventions, etc). Complex semantics, that span sequences of several assembly instructions, are hardly representable if embeddings of instructions are created in *isolation*; they have to be learned by the neural architecture using the embeddings.

Code is a means of communication between humans and machines, thus, unlike natural language, it has a dual nature. One level represents the syntactic and semantic meaning that can be inferred from its *static form*; the other level is the possibility of being *executed*. Some aspects of an Instruction Set Architecture (ISA) can be easily appreciated only when code is executed (e.g., the dependencies introduced by `RFLAGS` in X64). Moreover, semantically equivalent but syntactically different sequences of instructions are easily recognized when executed. Surprisingly, all embedding techniques we are aware of, only consider the static aspect of binary code.

1.2 Fine-tunable Assembly Model

A second limitation of current solutions stems from the fact that all proposed models are non fine-tunable (see use-case (a) of Figure 1). The DNN creates an embedding vector for each instruction, and this vector is then fed into the specialized neural architecture A that solves the target task. In this paradigm, there is no feedback from the specific problem to the embedding layer. All the semantically relevant information has to be stored into a single vector of real numbers. It is clear that there is a limit on the information encodable in a small fixed-size vector.

We propose a different approach based on a *fine-tunable* assembly model (use-case (b) of Figure 1). This is a transformer-based encoder [36] pre-trained on a large corpus of assembly code using specific tasks. During pre-training, the encoder learns a general semantic of assembly sequences that is context and execution aware. Then the pre-trained model is a part of a DNN A that solves a specific *downstream task* (e.g., compiler provenance, function similarity, and others). The DNN, including the assembly model, is retrained end-to-end on a small amount of problem-specific data during the fine-tuning process. Intuitively, during pre-training, the encoder stores the learned information in its internal weights (tens of millions against the few hundreds of an embedding vector), while during fine-tuning, weights are rearranged and changed to apply acquired knowledge for the final goal (i.e., *knowledge transfer*). This paradigm is state of the art for NLP and works well also if the fine-tune dataset is small. This is useful whenever the definition of a labeled dataset requires expensive manual effort.

1.3 Our proposal: BinBert

In this paper we introduce *BinBert*, a fine-tunable assembly code model based on a transformer encoder that is execution-aware. To inject execution awareness into our model, our idea is to symbolically execute snippets of assembly code. Specifically, we use a symbolic execution engine that transforms sequences of assembly instructions connected by a data-dependency relationship (the strands introduced in [12]) into sets of semantically equivalent symbolic expressions. These expressions are a functional representation of the input-output

relationship of the strand. We designed a novel pre-training process that forces BinBert to learn the correct matching between an assembly sequence and an equivalent symbolic expression and to translate assembly code into symbolic expressions and vice-versa.

We train BinBert on a new large dataset of assembly sequences and symbolic expressions derived from symbolic execution, obtaining a general-purpose assembly code model. We remark that symbolic execution is needed only in the pre-training phase; no code execution is required while using the model for inference tasks.

We tested BinBert on a multi-task benchmark for binary code understanding that we built. Tasks in the benchmark range from intrinsic ones, aimed at evaluating how the pre-trained BinBert captures the semantic of instructions and sequences, to extrinsic downstream tasks, in which we fine-tune BinBert for problems on assembly sequences and binary functions. In all our experiments BinBert raises the performance bar outperforming the current state of the art (including the recent PalmTree [19]) and specific solutions created for the binary similarity problem.

In summary, this paper provides the following contributions:

- a novel training task that makes the training of an assembly code model execution-aware by using symbolic expressions derived from the symbolic executions of assembly snippets;
- BinBert, a pre-trained execution-aware transformer model for X64, that can be plugged into DNNs for binary analysis;
- the first multi-task benchmark designed to test the binary code understanding of assembly models. The benchmark is composed of well-known tasks selected from the literature for their relevance, and two novel tasks (strand recovery and execution) for the semantic understanding of assembly sequences;
- an in-depth performance evaluation of BinBert based on our benchmark that shows how execution awareness improves the performance of an assembly model. To the best of our knowledge, we are the first to thoroughly test the impact of the fine-tuning paradigm on assembly. We show that, as already shown in the NLP field, the pre-training/fine-tuning approach has a huge positive impact on all downstream tasks. As a consequence, BinBert outperforms the current state-of-the-art instruction embedding technique.

2 Background

In this section, we introduce the general theoretical concepts behind the instruction embedding techniques and focus on the current state of the art. Afterwards, we detail weak points and gaps in current solutions, discussing how these influenced our proposal.

2.1 Instruction Embedding Models

An instruction embedding model takes as input an assembly instruction i from a vocabulary V of size d and it returns a vector of real numbers $e(i) = \vec{i} \in \mathbb{R}^n$, n is the embedding size (typically $n \in [128, 1024]$). The vector \vec{i} is a *dense representation* of the instruction i .

In the simplest embedding scheme a random matrix M of size $\mathbb{R}^{d \times n}$ is created, each instruction is mapped univocally to a row of M . With the matrix a sequence of instructions $I = [i_0, i_1, \dots, i_m]$ is converted into a sequence of vectors $e(I) = [\vec{i}_0, \vec{i}_1, \dots, \vec{i}_m]$ using a lookup mechanism. This sequence is fed into the task-specific DNN A . The matrix M is usually *trainable*: its elements are trainable weights and are modified during the training of A .

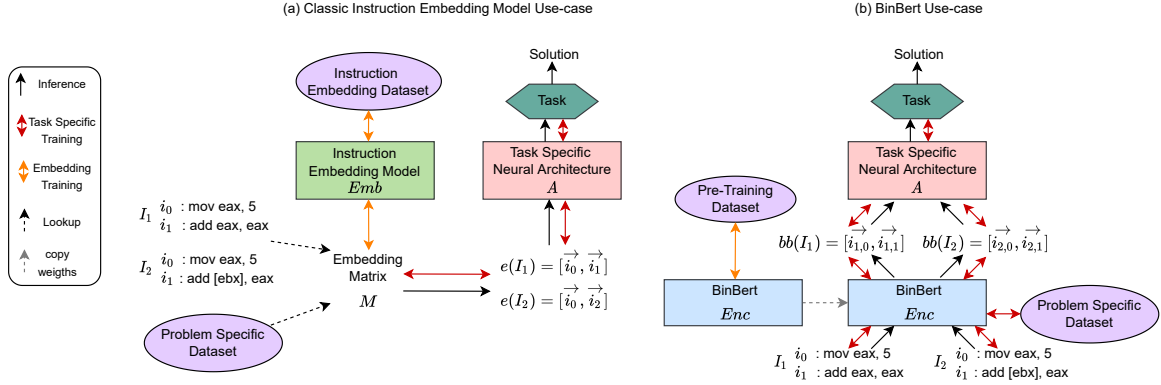


Figure 1: Comparison between an instruction embedding model (a) and the BinBert assembly model (b). In (a) the embedding model trained on a corpus C is a lookup mechanism that transforms instructions into vectors. Semantically relevant information have to be encoded in the embedding matrix M . The model cannot directly transfer knowledge into the final task. The BinBert assembly language model can use all the parameters of its neural architecture to store semantically relevant information that is then transferred to the downstream task with end-to-end training.

The groundbreaking idea of the embedding models is to generate the embedding matrix M with a neural network Emb , formally speaking $M = Emb(A)$. The network Emb is trained in an unsupervised way on a corpus C of data. The use case of the instruction embedding models is reported in Figure 1. This corpora C is composed of sequences of assembly instructions extracted from selected binaries. Usually, the *distributed representation learning tasks* used by instruction embedding models are, apart from minimal modifications, the ones used in NLP by solutions such as word2vec [24], GloVe [31], fastText [4], pv-dm [18]. The common goal is to train Emb to produce an embedding vector that contains enough information to predict a masked instruction from its context in C . Most of the novelty of the instruction embedding system is in the preprocessing of instructions and in the definition of the assembly sequences composing C .

2.1.1 Preprocessing of Assembly Instructions

Assembly language and natural language are distinguished by the wide difference between the vocabulary size d . A natural language is usually composed of hundred thousands different words, while the number of possible distinct assembly instructions is much more. Consider the X64 ISA, a mov instruction can use 64 bits to express immediates, offsets, and memory addresses, thus there can be 2^{64} different instructions that just move a value in a certain register. This makes raw assembly instructions impractical: a large vocabulary is discouraged [7] as it worsens the problem of out-of-vocabulary word (OOV) [17].

Moreover, the exact value of an immediate is largely useless in a static analysis setting (e.g. a memory address of an unknown memory layout) [22]. To ameliorate this problem, a lot of effort has been devoted to instructions preprocessing [42, 15, 30, 14, 19, 22]. The standard of the field is to substitute all memory addresses and immediates above a certain threshold value with special symbols (e.g. IMM). Another design choice is whether to consider the entire assembly instructions as a token (used in [22]), to split the assembly instructions into several tokens by separating opcodes and operands (used in [14]) or to use a more fine-grained split strategy [19]. Interestingly, no one used automatic tokenization such as WordPiece [39] that are standard

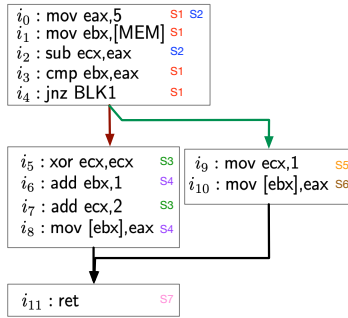


Figure 2: CFG of an imaginary function in x64 assembly. For each instruction we indicate with s_x the block level strand to which it belongs.

in NLP.

2.1.2 Extraction of Assembly Sequences

A key point is how to extract the sequences from the binary, as this defines the context in which an instruction appears. The context for instruction i_x is composed by k instructions appearing before/after i_x in C . Used extraction strategies are:

- **Linearized CFG:** In this case each sequence is a linearization of a CFG (commonly the one provided by a disassembler) [22]. Blocks of the CFG that are not logically related could be sequentially placed in the linearization, and thus an instruction will see a noisy context. An example is the linearization of the CFG in Figure 2 induced by instructions numbers: the context of instruction i_8 contains instructions i_9, i_{10} that are not causally related. This inject noise in the learning process.
- **CFG/ICGF:** the sequences are extracted from the recovered CFG/ICFG. This is done either by using a random walk strategy ([14]), or by taking as a sequence a single block [42]. The idea is to have a sequence that respects the logical control flow of the examined program, removing the source of noise highlighted in the previous strategy. We argue that this technique does not completely remove the presence of extraneous instructions in the context. Take the sequence of instructions i_5, i_6, i_7, i_8 in Figure 2, the context of instruction i_7 contains i_6 and i_8 that are not causally related.

2.1.3 PalmTree

PalmTree [19] is an instruction embedding model that has shown state-of-the-art performances beating all the other embedding models on several tasks. PalmTree is based on a transformer model. PalmTree divides an instruction into tokens using a fine-grained strategy with manually made regexs. PalmTree is trained on pairs of instructions taken from the corpora C . PalmTree uses the standard MLM of Bert and two novel tasks: the CWP in which the network has to recognise if a pair of instructions is taken from the same context or not, and the DUP in which the network has to recognize if there is a data dependency between instructions. Once trained, PalmTree is used as an instruction embedding model (see use-case (a) in Figure 1): a sequence of instructions is embedded by applying separately the PalmTree model to each instruction.

2.2 Weak Points and Gap Analysis

We can now identify weak points of previous solutions and gaps in the current body of knowledge. From such analysis, we derive the research directions of our solution.

2.2.1 Focus on a Single Assembly Instruction and No Context-Awareness

Current solutions transform a single assembly instruction into a vector. This approach, despite its proven effectiveness [9, 22], has an implicit weakness: the semantic of a single assembly instruction is really limited. Even a perfect embedding model is constrained to produce embedding vectors that capture only the semantic of a single instruction. However, non-trivial concepts are encoded using an entire sequence of instructions (e.g., swap of two variables in memory, spilling/restoring of variables in/from the stack,...).

On the other hand, the meaning of an assembly instruction depends from its context. Consider the sequences I_0 : `lea eax, [ebx * ecx + edx]; mov edi, eax` and I_1 : `lea eax, [ebx * ecx + edx]; mov edi, [eax]`; in the first case the `lea` is used to perform arithmetic operations ($edi = (ebx * ecx) + edx$), in the second case is used to compute a memory address. An instruction embedding approach cannot encode these concepts, even if they are present in the corpora C . Therefore, (re)learning such patterns is demanded to the upward neural architecture A (see Figure 1). This means that A has to be trained on a large enough problem-specific dataset, even if Emb was trained on large corpora of data. Thus our first research direction is:

RD 1: Design an embedding model that could transform entire sequences of instructions into semantically representative embedding vectors. If the model is used to generate the embedding of an instruction, it has to take into account its context.

2.2.2 No Execution

The execution of assembly makes clear concepts that would be covert (or unavailable) from its static representation. Apart from the `RFLAGS` example mentioned in the introduction, consider again the sequence I_0 : `lea eax, [ebx * ecx + edx]; mov edi, eax` such sequence is semantically equivalent to I'_0 : `imul ebx, ecx; add ebx, edx; mov edi, eax`. The semantic equivalence could be easily discovered if information taken from the execution is inserted in the pre-training tasks. Unfortunately, current techniques totally neglect this aspect.

RD 2: Design an embedding model that takes into account the execution of code. The impact of execution-related information has to be quantified with a specific ablation study.

2.2.3 No End-To-End Retraining (Fine-tuning)

The use case of all the known instruction embedding models is the (a) of Figure 1. In the NLP field, such approach has been largely substituted by encoder models pre-trained on corpora C and then fine-tuned on a specific task (use-case (b) of Figure 1). During fine-tuning, the pre-trained encoder Enc is trained with the network A end-to-end on a problem specific labeled dataset. The back-propagation algorithm optimizes the internal weights of network A and Enc ; this optimization modifies the weights of Enc so that the knowledge learned during pre-training is applied to the downstream task. As a matter of fact, when Enc is a transformer, A is usually a linear classifier or another simple neural network, since most of the work is performed by Enc .

No one has extensively studied how this paradigm copes with solving different goals on assembly language, goals selected to test the semantic and syntactic comprehension of assembly code. This interesting gap of the current body of knowledge gives us a new research direction:

RD 3: Design an embedding model that can be trained end-to-end on a specific task, transferring the general knowledge learned during pre-training. The model has to be evaluated on a multi-task benchmark designed to thoroughly test the syntactic and semantic understanding of the assembly language.

3 The BinBert Solution

In this section we describe BinBert. We first give an overview of the system, briefly describing the transformer architecture [36]. We then give the details of the innovative aspects of BinBert.

3.1 Overview

The neural architecture of BinBert is the standard transformer encoder [36] used by Bert [13]. A transformer encoder processes sequential data using an attention mechanism, which allows for both the creation of more informative embeddings (by focusing only on relevant parts of the sequence) and good performances (the attention mechanism is implemented using matrix multiplication that is highly parallelizable on GPUs). In detail, a transformer encoder is composed of N identical layers stacked one on top of the other, where each layer consists of two sublayers: a multi-head self-attention mechanism and a fully connected feed-forward network. Practically speaking, a sequence of n tokens is transformed into a sequence of n latent vectors (one for each token) with a mechanism that we will explain in Section 3.4; this sequence is fed into the initial layer of the encoder. Each other layer takes as input the hidden state token vectors returned by the previous layer. The output of the encoder is a sequence of $n + 1$ embedding vectors: one for each input token and a special embedding for the entire sequence (the [CLS] vector described in Section 3.4).

To avoid the vocabulary inflation generated by the use of raw assembly instructions, BinBert substitutes memory addresses and immediates above a certain threshold with special symbols. Moreover, BinBert splits a single assembly instruction into several tokens using WordPiece [39].

In BinBert we decided to completely remove the noise given by instructions that are contextually related but have no logical relation (see Section 2.1.2) by extracting sequences representing *strands* [12]. Strands are sequences of causally related instructions computing the values of a certain variable. In this way, the context of an instruction never contains extraneous instructions introduced by compiler optimizations. We symbolically execute each strand to extract a set of symbolic expressions; these expressions will be used in our training tasks as a means to inject execution-related information into the pre-training.

During pre-training, BinBert learns the matching between symbolic expressions and strands (this is done using positive/negative pairs); at the same time, samples are partially masked according to a translation task.

3.2 Instructions Preprocessing and Assembly Sequences Extraction

We preprocess each assembly instruction substituting immediates above a threshold (5000 in our experiments) with the value IMM (the same is done for offsets and memory addresses). We use the special symbol MEM in case of jumps. We use a threshold-based approach as small immediate values are like to carry informative content (comparison with small constants in if and loops, PC/stack relative displacements that identify variables in memory). All immediates/offsets are converted to decimal format. For call instructions, we distinguish if the called function is user-defined or belongs to libc. For user-defined functions, we substitute the called address

with `func` (our system is usable on stripped binaries). If it is a call to `libc`, we substitute the address with the function name (e.g., `call printf`), since external symbols cannot be stripped. Indirect calls are left untouched.

After preprocessing, each instruction is tokenized using WordPiece [39]. The latter uses a probabilistic approach to learn how to tokenize instructions in a way that minimizes the vocabulary size and the OOW problem. Contrarily to manually made regexes, WordPiece automatically learns how to split complex opcodes (as an example `cmovz` will be split in `cmov` and `z` helping the model in understanding the relationship between the `cmovX` family of X64). We use WordPiece also on symbolic expressions, This provides a uniform tokenization mechanism and vocabulary for the two distinct languages (`asm/sym. expr.`) used for BinBert.

Assembly Sequences - Strands extraction In BinBert we use the concept of strands to extract the sequences of assembly instructions on which our model is trained. This does not mean that BinBert cannot be fine-tuned and used on CFG blocks or entire functions as our experiments will show.

A strand, originally defined in [12], is a slice of a CFG block constituted by all the instructions that are connected by def-use dependences. More specifically, we consider as an output variable of a block a memory location or a register on which the last operation is a write or the check of a jump. Starting from this variable we construct a *backward-slice* of the block including all the instructions from which the value of such variables depends. To make the concept clear, consider the example in Figure 2; the first block contains two strands S1 and S2: S1 is composed by instructions i_0, i_1, i_3, i_4 that influence the `RFLAGS` register later checked by i_4 ; strand S2 is composed by i_0, i_2 which define the value of variable `ecx`. Other examples of strands are in the figure. We enrich the original definition of strand, by considering as a single strand all the instructions that prepare the input values for a call. In this case, the strand will be constituted by all the aforementioned instructions and the call instruction itself. Therefore, we build our training corpora C by extracting the CFGs in binary, and then decomposing all the blocks in strands. The strands will be the basic sequences in C .

This decomposition has several advantages: it completely removes the noise introduced by instructions that co-occur in the same context only for compiler optimization reasons and the model learns the entire “*causal context*” of an instruction so it is able to see long dependencies among instructions.

3.3 Symbolic Execution

In BinBert, we use symbolic execution to convert each strand into a representative symbolic expression. The symbolic execution engine works on strands of assembly instructions before we apply the preprocessing (it needs the actual value of immediates and memory locations). The engine is built on `angr` [34]. During the execution we consider as inputs the variables on which the first operation executed by the strand is a read, and as outputs the variables on which the last operation executed is a write.

Each time the strand writes on a variable (either a memory location or a register) we express the written value using a symbolic expression. When a variable is read we may have that either the variable is an input (no one wrote on it) or it contains a symbolic value. In case the variable is an input, we set its symbolic value to its address or register name (as example, for `mov eax, [rbp+4]` we have `eax=*(rbp+4)`, for `mov eax, ebx` we have `eax=ebx`).

The symbolic expression obtained with this process may have one of three possible forms:

- If the strand computes the value of a certain variable, the symbolic expression describes the value in the output variable as a function of the strand inputs. An example is in the first row of Table 1: we have the symbolic expression `rcx=-1 add (0 Concat rsi[1:0])`. This expression is for the output variable `ecx`: the `and` operation extracts the two least significant bits from `rsi`, the result is extended

Strands	Symbolic Expressions
mov ecx, esi	rcx = -1 add (0 Concat rsi[1:0])
and ecx, 3	rdi = 1 add rdi
rep stosb byte ptr [rdi], al	*(rdi) = al
mov rax, qword ptr [rbp - 168]	0 Concat (*(rbp add -168) add 24)[1:1] ne 0
mov eax, dword ptr [rax + 24]	
test al, 2	
jne MEM	
mov esi, IMM	fprintf(*(rbp), IMM)
mov rdi, qword ptr [rbp]	
call fprintf	

Table 1: Examples of symbolic expressions obtained by different strands.

to 64 bits, and then decremented by the `rep`. The expression only contains the extended register of X64, this is a design choice that we discuss later in Section 3.4.

- If the strand computes the predicate checked by a conditional branch, then our symbolic expression will be the comparison of the jump condition with a symbolic expression of the value used in the predicate. As an example, consider the second row of Table 1; in this case, the symbolic expression is compared with 0 using the not equal (ne) predicate.
- Finally, in case our strand computes the arguments used by a call instruction, our symbolic expression will be the call to the specific function (including the symbolic name if it is a libc call), where all the arguments are substituted by the symbolic expressions of their values. We extracted function arguments following the X64 calling convention. An example of a call expression is in the third row of Table 1.

From a single strand, we may have several symbolic expressions. For instance, in the first row of Table 1 the `rep stosb` instruction, repeatedly puts the content of `al` into the memory pointed by `rdi` for `ecx` times by decrementing `ecx` and incrementing `rdi` at each iteration. This means that the strand has three output variables (a memory location and two registers), and each one will have its symbolic expression. We will call such set of symbolic expressions the *representative set* of the strand.

Preprocessing and Tokenization of Symbolic Expressions The symbolic expressions of each strand are preprocessed similarly to assembly. Large numerical constants are substituted with the special symbol `IMM`, for floating point numbers all the digits after two decimals are truncated. The symbolic expressions are tokenized using `WordPiece`. During the preprocessing phase, we also substitute the name of all registers to their extended form (i.e., we use `rax` instead of `eax`), we do so to help the network in understanding the relationships between names used to address different parts of the same logical register; this step is not applied while preprocessing assembly instructions (i.e., we leave `eax` in the strand).

3.4 BinBert Input Representation and Pre-Training Tasks

BinBert is fed with pairs `<strand, symbolic expression>` and pre-trained on two tasks that we name Execution Language Modeling (ELM) and Strand-Symbolic Mapping (SSM).

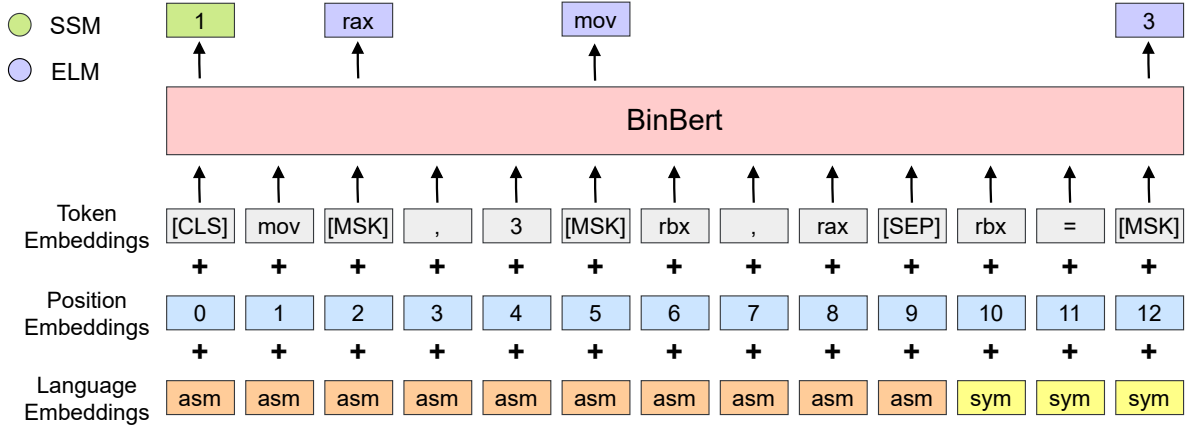


Figure 3: Example of BinBERT pre-trained on the Execution Language Modeling (ELM) and on the Strand-Symbolic Mapping (SSM) tasks.

Input Representation A BinBERT input is a tokenized strand-symbolic expression pair. Two special tokens are added to each sample: [SEP] used to keep the distinction between assembly and symbolic expression, [CLS] prepended to all the samples. The hidden state of the [CLS] token in the last hidden layer is usually used to obtain a latent vector representation of the whole sequence [13] (see Figure 3).

We use dynamic padding, thus shorter sequences are padded with the special token [PAD] and sequences longer than a threshold (we use 512 in our experiments) are truncated. Prior to the transformer architecture, each token is converted into a vector by using the lookup mechanism described in Section 2.1. Then, token embeddings are summed up with both position and language embeddings [13, 10]. Position embeddings are used to make the model aware of the notion of sequences, while language embeddings are used to distinguish assembly from symbolic expressions (see Figure 3).

Pre-Training Tasks The first task is **Execution Language Modeling (ELM)**. The goal of ELM is to mask a certain percentage m_p of tokens in the input pairs and to make the network predict the original ones. As in Bert, the tokens that have to be predicted are either substituted with a special token [MASK], replaced with a random word, or left untouched with 80-10-10 probabilities respectively. Figure 3, shows an example in which the tokens `rax`, `mov`, and `3` are masked and the network tries to reconstruct them as output. Notice that, in order to reconstruct a token contained in a strand, the network needs to pay attention to both the strand and its corresponding symbolic expression (the same holds for a token in the symbolic expression). This means that in order to predict the value `3` for register `rbx` the network is forced to understand the data flow in the strand, thus making the semantic of each strand instruction explicit. As in [13] a linear layer is added on top of the last hidden layer, in particular to the hidden states of the masked tokens; this layer will guess the masked tokens.

Mathematically speaking, we have a dataset \mathcal{D} of $\langle \text{strand}, \text{symbolic expression} \rangle$ pairs $I = [t_1, \dots, t_n] \in \mathcal{D}$ where each token t_i belongs to a vocabulary $V \in \mathbb{R}^d$. We feed each pair to BinBERT to obtain context-aware hidden vectors $H = \text{Binbert}(I) = [\vec{h}_1, \dots, \vec{h}_n]$ as output. Now, consider a function $\text{msk}(I, m_p)$ which randomly mask m_p percentage of tokens of I , then the goal of the network is to predict the probability that a

token t_i corresponds to the target word \hat{t}_i with the softmax function:

$$p(t_i = \hat{t}_i | I) = \frac{e^{\vec{w}_{\hat{t}_i} \cdot \vec{h}_i}}{\sum_{k=1}^d e^{\vec{w}_{\hat{t}_k} \cdot \vec{h}_i}} \quad (1)$$

where $\vec{w}_{\hat{t}_i}$ is the weight vector of the linear layer for word \hat{t}_i . The loss of the ELM task is the cross entropy loss:

$$\mathcal{L}_{ELM} = - \sum_{I \in \mathcal{D}} \sum_{t_i \in \text{msk}(I, m_p)} \log p(t_i = \hat{t}_i | I) \quad (2)$$

The second task is the **Strand-Symbolic Mapping** (SSM), in which, the goal is to predict whether the symbolic expression is from the set of expressions representative of the strand (see Section 3.3). To solve this task, we create both negative pairs by associating a strand with a random symbolic expression and positive pairs in which the symbolic expression is taken from its representative set. The ratio between positive and negative pairs is 50:50. An example of positive pair can be seen in Figure 3: the input strand computes the value 3 for register *rbx* as stated by the corresponding symbolic expression. We believe that, with this task, the network is forced to learn the matching assembly snippets and symbolic expressions. This task is built by using a linear layer on top of the hidden state of the [CLS] token in the last layer that will be used to classify a pair as negative or positive. In mathematical terms, the goal of this task is to evaluate the probability that the output label is one, i.e. the symbolic expression correctly computes a value in the strand:

$$p(y = 1 | I) = \frac{e^{\vec{w}_1 \cdot \vec{h}_{[CLS]}}}{e^{\vec{w}_0 \cdot \vec{h}_{[CLS]}} + e^{\vec{w}_1 \cdot \vec{h}_{[CLS]}}} \quad (3)$$

where \vec{w}_0 and \vec{w}_1 are the weight vector of the linear layer for label 0 (negative pair) and 1 (positive pair) respectively. The loss \mathcal{L}_{SSM} of the SSM task is the standard cross entropy loss.

The final loss on which BinBert is trained is the sum of the losses of the two tasks described above:

$$\mathcal{L} = \mathcal{L}_{ELM} + \mathcal{L}_{SSM} \quad (4)$$

4 Evaluation Tasks

When proposing a new assembly model an extensive experimental evaluation is fundamental. In NLP is customary to use standard multi-task benchmarks to evaluate language models [37]. Similar benchmarks for binary code do not exist. Therefore, we designed our benchmark by selecting several tasks. For each task defined on a sequence of assembly instructions, we have a version on strands and one on CFGs' blocks. This tests BinBert on sequences that are not strands.

4.1 Intrinsic Tasks

The intrinsic tasks [6] directly use the embeddings produced by BinBert, there is no fine-tuning and the embeddings are not used as input for other models.

4.1.1 Opcode and Operand Outliers

In the **opcode outlier** task we are given a set of five instructions. Four of these instructions belong to the same semantic class, and one is an outlier. As example, if the set is `{add eax, ebx; sub ebx, ecx; imul ecx, edx; add eax, 5; call printf;}`, the last one is an outlier. The network has to predict which is the outlier, among the 5 instructions.

The **operand outlier** task is analogous, but in this case, the operands define the outlier. Given the set `{sub [eax+5], ebx; sub ebx, ecx; imul ecx, edx; add eax, ebx, add ebp, esp;}`, the outlier is the first instruction: the only using a memory operand. Again, we followed the operand categorization provided by [19].

4.1.2 Strand and Block Similarity

In the **strand similarity** task we compute the embeddings of given strands with BinBert and use them to discover semantically similar strands. Two strands are similar if they have an overlapping semantic (non-empty intersection of the representative sets). More formally, we have a lookup database of n strands $\mathbb{A} = \{a_1, \dots, a_n\}$ and a query strand q . The lookup contains strands that are similar to q and strands that are dissimilar. Given a number k the network has to return the k strands in \mathbb{A} that are most similar to q . The **block similarity** is analogous but done at basic block level. In this case, we use the definition of similar blocks used in [15], where two blocks are similar if the DWARF information says so.

4.2 Extrinsic Tasks

In the extrinsic tasks, BinBert will be used as the encoding layer of a neural architecture and fine-tuned end-to-end.

4.2.1 Strand and Block Similarity

The **extrinsic strand similarity** and **extrinsic block similarity** are the same tasks of their intrinsic versions; in this case, we fine-tune BinBert using a dataset of similar and dissimilar pairs of samples. We decide to include these tasks as they will quantify the effect of fine-tuning on the creation of semantic preserving embeddings. These tasks are the dual of the compiler provenance below.

4.2.2 Strand and Block Compiler Provenance

In the **strand compiler provenance** task the architecture has to recognise the compiler and the optimization levels used to generate a particular strand. This task has been previously proposed on functions [23] and fragments of code [25]. In compiler provenance, the network has to recognize the syntactic signature that a compiler, or optimization level, produces while in the strand and block similarity it has to abstract from such differences to get the meaning of the sequences. The **block compiler provenance** task is analogous but at block level.

4.2.3 Strand Recovery and Execution

We designed two novel tasks that test the semantic understanding of assembly sequences. In **strand recovery** we provide to the network a basic block of the CFG where one instruction is marked. The DNN has to recognize all instructions in the same strand of the marked instruction. This task tests the understanding of the inputs/outputs of instructions, the network has to infer the dependency created by implicit registers such as RFLAGS.

In **strand execution** a strand and a question are given to the network. The question is composed of an assignment for the inputs and a marked output. The network has to predict the value of the marked output. This task is interesting as it forces the network to concretely execute the snippet of assembly code and compute the correct output.

4.2.4 Function Level - Compiler Provenance and Similarity

Finally, our multitask benchmark contains two tasks at the function level. In the **function compiler provenance** task the network is given an entire binary function and it has to predict the compiler used to generate the function and the optimization level. The **function similarity** task is analogous to the extrinsic block similarity: given a database of functions and a set of function queries, for each query, the network has to return the similar functions in the dataset. We use the standard definition of function similarity [22]: two assembly functions are similar when they derive from the same source code compiled with different compilers or optimization levels. Function similarity is a current hot topic in research [22, 42, 15, 41, 30] for its security implications.

5 Datasets, Pre-Training and Implementation Details

5.0.1 Datasets

We used two datasets, a pre-train dataset PTData used to pre-train BinBert, and a test dataset TestData for the downstream tasks. Our datasets are for X64.

PTData - Pre-train Dataset The dataset contains the projects: ccv-0.7, binutils-2.30, valgrind-3.13.0, libhttpd-2.0, openssl-1.1.1-pre8, openmpi-3.1.1, coreutils-8.29, gsl-2.5, gdb-8.2, postgresql-10.4, ffmpeg-4.0.2, curl-7.61.0. The projects are compiled using: clang-3.8, clang-3.9, clang-4.0, clang-5.0, gcc-3.4, gcc-4.7, gcc-4.8, gcc-4.9, gcc-5.0. For each compiler we compiled each project four times, one for each optimization level in $\{O0, O1, O2, O3\}$. We use radare2 (version 5.6.0) to extract functions signatures. We use angr (version 9.1.11611) to get CFGs, basic blocks, and strands and to obtain from each strand the set of symbolic expressions. After the removal of duplicates, we obtained 17.215.046 pairs in the form $(strand, simexpr)$ as dataset for the SSM task.

TestData - Test Dataset The test dataset is obtained from diffutils-3.7, findutils-4.7.0, inetutils-2.0, mailutils-3.10, wget-1.20.3. We use clang-3.8, clang-6.0, clang-9, gcc-5, gcc-7, gcc-9, icc-21 and the 4 optimization levels $\{O0, O1, O2, O3\}$ to obtain the raw binaries. We will use these raw binaries to create the specific dataset for each task. Since some operations are task-dependent we discuss the specific split and format of the data in the experimental section of each task. We took care of removing duplicates so that the same sample will not be in fine-tune and test split.

5.0.2 Model Parameters, Pre-Training and Implementation Details

Our model is built using python 3, with pytorch (version 1.10.2+cu113) [27] and huggingface (version 4.16.02) [38]. We trained it on a DGX A100, using 4 A100 GPUs.

Model and Pre-Training Parameters BinBert parameters are: sequence length 512, hidden size 768, intermediate size 3072, 12 attention heads and layers. We used AdamOptimizer and learning rate 0.0001. The

masking rate mp is 0.3. The batch size for each device is 32 with two steps of gradient accumulation; having 4 GPUs the equivalent batch size is 256. We trained for 1 epoch using 1425 steps of warmup.

6 Experimental Evaluation

In our evaluation we answer the following experimental questions:

RQ 1 Is an execution-aware transformer model trained on strands of assembly instructions the state-of-the-art assembly model for binary understanding?

RQ 2 What is the impact of pre-train on the performance across several binary understanding downstream tasks?

RQ 3 What is the impact of using an execution-aware pre-training?

To answer **RQ 1** we compare BinBert with PalmTree on all the tasks of our benchmark. We also compare BinBert with state-of-the-art function similarity solutions Safe [22] and asm2vec [14] on the extrinsic function similarity task. For a fair comparison, we retrain both PalmTree and asm2vec on our pre-train dataset, by using the same parameters of the original papers. On all the extrinsic tasks we fine-tune PalmTree¹ on exactly the same dataset we used to fine-tune BinBert. We are interested in assessing the effective contribution of the symbolic expressions used during pre-training (**RQ 3**) and the pre-training itself on downstream applications (**RQ 2**). To do so we will use the following baselines:

- **BinBert-MLM**: it is a model pre-trained on strands only with the standard Masked Language Modeling Task (MLM). In MLM only the assembly of a strand is given to the transformer during pre-training, the pre-training task is to recover masked tokens. The gap between this model and BinBert quantifies the impact of execution-awareness.
- **BinBert-FS**: it is a transformer encoder with the application-specific neural architecture trained from scratch on the specific downstream task. The gap between this model and BinBert quantifies the impact of pre-training.
- **PalmTree-FS**: similar to BinBert-FS but for PalmTree. The gap between this model and PalmTree quantifies the impact of pre-training.

For the fine-tuned models we use the notation: **BinBert-FT**, **BinBert-MLM-FT**, and **PalmTree-FT**.

For each task, we will provide the details of the dataset, the solution we employed to solve the problem, the metrics used to evaluate the performance, and the final results.

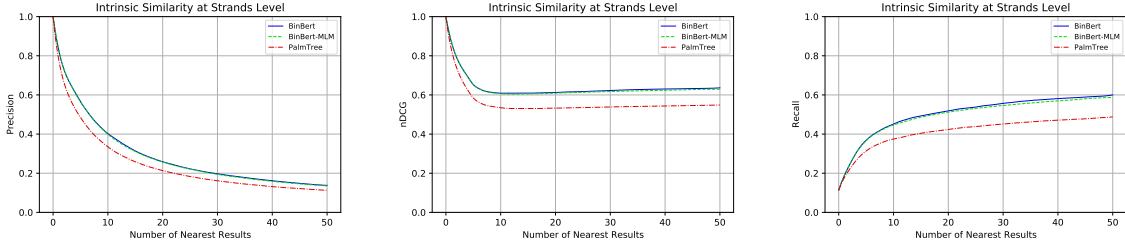
6.1 Intrinsic Tasks

The results for intrinsic tasks are reported in Table 2.

¹PalmTree authors highlighted the possibility of fine-tuning their system but they did not explore this possibility in their paper.

			Models								
			BinBert			BinBert-MLM			PalmTree		
			Accuracy			Accuracy			Accuracy		
Outlier Detection	Opcode		83.5 ± 0.17			72.6 ± 0.12			79.2 ± 0.25		
	Operand		86.2 ± 0.17			74.8 ± 0.19			71.1 ± 0.28		
			Prec.	Rec.	nDCG	Prec.	Rec.	nDCG	Prec.	Rec.	nDCG
Similarity	Strands	top-10	42.7	44.0	61.1	42.4	43.6	60.9	35.7	36.8	53.9
		top-25	22.9	53.5	61.6	22.7	52.8	61.3	18.9	43.6	53.5
		top-40	16.5	57.9	63.0	16.1	56.7	62.3	13.4	47.0	54.3
	Basic Blocks	top-5	55.3	25.4	61.2	54.9	25.0	61.3	52.7	24.1	59.2
		top-10	38.9	34.5	50.1	38.3	33.9	49.9	35.8	31.8	47.4
		top-25	19.5	42.5	48.9	19.5	42.2	48.9	18.2	39.6	46.4

Table 2: Intrinsic evaluation results.



(a) Precision for the top- k answers with $k \leq 50$. (b) nDCG for the top- k answers with $k \leq 50$. (c) Recall for the top- k answers with $k \leq 50$.

Figure 4: Results for the **intrinsic strand similarity** task. Database of 49079 strands, average on 833 queries.

6.1.1 Opcode and Operand Outlier

We used 43879 different instructions to create a dataset of 50000 sets of 5 instructions. These sets are created according to the task (either opcode or operand) as we have defined in Section 4.1.1. We use opcode and operand categories analogous to [19]: opcodes are categorized according to the *x86 Assembly Language Reference Manual*², while operands are categorized according to their type (e.g. two registers, one register and one memory access, etc.). Further details about these classes can be found in the Appendix 9.1. To solve this task, we embed each instruction in the set and we evaluate if the embeddings are able to distinguish the outlier; this is done by computing the distance of each embedded instruction from the others and by predicting as outlier the most distant. An instruction embedding is computed by mean pooling the instruction tokens’ hidden states in the second last layer of BinBert (we take the second last layer as it is less influenced by the pre-training task). The evaluation metric is the *accuracy* of [6]:

$$Accuracy = \frac{\sum_{s \in \mathcal{S}} outlier(s)}{|\mathcal{S}|} \quad (5)$$

²https://docs.oracle.com/cd/E26502_01/html/E28388/ennbz.html

where \mathbb{S} is the dataset composed of instruction sets and $outlier(s)$ is equal to 1 if the outlier in the instruction set s is detected and 0 otherwise.

We compute the mean accuracy and standard deviation on 10 runs of the experiment on different datasets (each composed of 50k sets), the results are in Table 2. BinBert achieves the best performances (0.84 on opcodes/0.86 on operands), it shows a great improvement over BinBert-MLM (0.73 on opcodes/0.75 on operands) this confirms that symbolic expressions clearly enrich the semantic learned by the model for each instruction. It also outperforms PalmTree (0.79 on opcodes/0.71 on operands) by a wide margin; the performances of PalmTree are not dominated by a transformer trained on MLM. We believe that this is so because PalmTree has been explicitly designed to be an instruction embedding solution, while BinBert-MLM has been trained on sequences. In Appendix 9.1 we report a qualitative analysis with the clusters of opcodes learned by BinBert.

6.1.2 Similarity at Strand Level

We use a database \mathbb{A} of 49079 strands, from this database we extract a database of 833 queries \mathbb{Q} . On average for each query we have 59 similar elements in \mathbb{A} . To solve the task we compute an embedding vector \vec{q} for each query q and an embedding vector \vec{a} for each strand a in the lookup database \mathbb{A} . This is done by averaging all the instruction tokens’ hidden states in the second last layer of BinBert. For each query vector \vec{q} we compute the cosine similarity with all $\vec{a} \in \mathbb{A}$; we return the ordered list of the top- k similar elements $R_q = (r_1, \dots, r_k)$. Using R_q we compute: precision, number of true similars in R_q over k ; recall, number of true similar in R_q over $\#sim(q)$, that is the number of items similar to q in \mathbb{A} ; and nDCG. The nDCG is a measure used in information retrieval. It is defined as:

$$nDCG = \frac{\sum_{i=1}^k \frac{sim(r_i, q)}{\log(1+i)}}{\sum_{i=1}^{\#sim(q)} \frac{1}{\log(1+i)}} \tag{6}$$

Where $sim(r_i, q)$ is 1 if q is similar to r_i and 0 otherwise. The quantity at the denominator is the scoring of a perfect answer, and the number at the numerator is the scoring of our system. The nDCG is between 0 and 1, and it takes into account the ordering of the items in R_q , giving better scores when similar items are ordered first. As an example let us suppose we have two results for the same query: $(1, 1, 0, 0)$ and $(0, 0, 1, 1)$ (where 1 means that the corresponding index in the result list is occupied by a similar item and 0 otherwise). These results have the same precision (i.e., $\frac{1}{2}$), but nDCG scores the first better. We average the per-query precision, recall, and nDCG to obtain the final metrics. Results are shown in Table 2 and in Figure 4.

We can see that BinBert achieves the best performance on precision, recall and nDCG.

6.1.3 Similarity at Block Level

To solve this task we used a database \mathbb{A} of 6460 basic blocks where 505 are queries \mathbb{Q} . On average each query has 10 similar elements in \mathbb{A} . The test procedure and the metrics are the same of the strand similarity case except for the basic block embedding computation. A basic block is first decomposed into strands and then its embedding is obtained as the average of its strand embeddings. Results are shown in Table 2 (the figures for the values of $k \in [0, 30]$ are reported in the Appendix). We can see the same general behavior observed with the strand similarity, with BinBert being the top performer. But there is less gap with respect to PalmTree. We believe that this is due to the fact that PalmTree is not trained directly on strands, and so the strand similarity task is harder for PalmTree than block similarity. It is interesting to note that the reverse effect is not observable for BinBert, the performances on blocks are slightly better than the ones on strands.

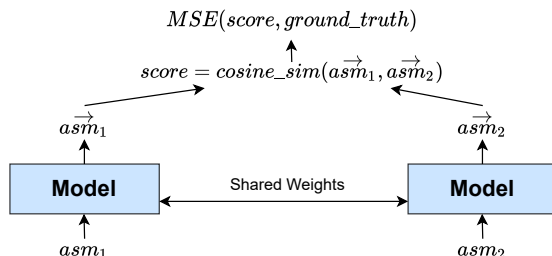


Figure 5: Scheme for the siamese architecture: asm_i can be either a strand, basic block or a function and Model is either BinBert, BinBert-ML or PalmTree.

Our evaluation shows that the execution-aware transformer BinBert is state-of-the-art for all the intrinsic tasks of our benchmark. Execution-awareness improves performances on all tasks. The benefit is greater on tasks that measure the understanding of instructions semantic.

6.2 Extrinsic Tasks at Strand and Block Level

Results for extrinsic tasks are shown in Table 3.

6.2.1 Strands and Blocks Similarity

Fine-tune Process and Dataset We fine-tuned BinBert, BinBert-MLM, and PalmTree on the task of recognizing if a pair of entities is similar or dissimilar. Considering strands, we construct a dataset of 44978 pairs of strands, 50% similar and 50% dissimilar. We split it into train and validation, resulting in 39979 strands pairs for the training set and 4999 for the validation. For blocks, we construct an analogous dataset with 24001 basic blocks pairs for the training set and 3001 pairs for the validation.

We fine-tune using a siamese architecture [5, 22, 40] (see Figure 5). In this architecture two instances of the embedding network are used, each instance produces the embedding of the corresponding entity in a pair. The resulting embeddings are compared with the cosine similarity to produce a score. The siamese network is trained with the *Mean Squared Error Loss (MSE)* which minimizes the distance between the similarity score produced by the network and the ground truth label that can be in $\{+1, -1\}$ (i.e. two strands/basic blocks are similar or not). Intuitively, this training process instructs the embedding network to produce embeddings that are close if they are from a similar pair and distant if not. For BinBert-FT and BinBert-MLM-FT the embedding that enters the cosine similarity is the average of all tokens of the last layer (excluded the padding).

Since PalmTree is an instruction embedding model, an additional architecture is needed to transform it into a model that embeds sequences. In particular, we use a bidirectional LSTM, where each cell takes as input the instruction embedding produced by PalmTree. Finally, we compute an embedding by averaging all the hidden states of the LSTM. We test the fine-tuned models on the same test sets used in the corresponding intrinsic tasks (see Sec.6.1). We fine-tune each model for 20 epochs, selecting the epoch with the best AUC on validation.

Results Results for strands and basic blocks similarities are in Table 3. In Figures 6 there are the results of strand similarity for $k \in [0, 50]$, the figures for block similarity are in the Appendix.

Also in this case BinBert achieves the best performances. The gap between BinBert-FT and the other models is much wider than in the intrinsic case. A possibility is that BinBert learns a wider semantic during pre-training,

solving more efficiently the similarity task after fine-tuning. This explains the gap between BinBert-FT and BinBert-MLM-FT.

The great impact of pre-training can be appreciated by looking at the performance of BinBert-FT and BinBert-FS. BinBert-FS has been only trained on the fine-tune dataset so it cannot leverage a learned semantic, the fine-tune dataset has not enough data to make up for this disadvantage and to train a big transformer model. The difference in model size between PalmTree-FS and BinBert-FS is the reason why PalmTree-FS performs better than BinBert-FS, a smaller model can be trained with less data.

The impact of pre-training can also be seen on PalmTree, PalmTree-FT beats PalmTree-FS. In this case, the difference is less marked than the one between BinBert-FT and BinBert-FS, the reasons probably are: PalmTree is a smaller model (and so it can be trained and reach a plateau with less data), and being an instruction embedding model cannot learn during pre-training the many complex interactions between instructions that can be useful for a sequence similarity task.

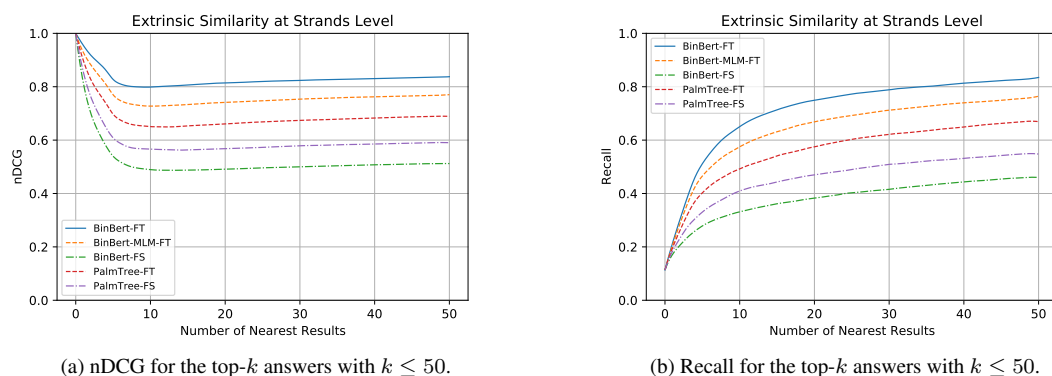


Figure 6: Results for the **extrinsic strand similarity** task. Database of 49079 strands, average on 833 queries.

6.2.2 Compiler Provenance

As in other works on compiler provenance [23], we train and test the networks on the task of *Compiler Classification*, that is detecting the compiler family that has generated a sample, and *Optimization Classification*, detecting the optimization level used to generate a sample.

Fine-tune Process and Dataset The compiler provenance dataset is made up of samples in which, depending on the granularity, a strand or a basic block is associated with two labels: the compiler family with its corresponding version and the optimization used. The strand compiler provenance dataset has 36828 samples, split into 29462 samples for the training set and 3683 samples for both the validation and test sets. The basic blocks compiler provenance dataset has 81918 samples, split with the same ratio as before, resulting in 65534 training samples and 8192 validation and test samples.

We fine-tune BinBert and BinBert-MLM on both compiler and optimization classification by adding a linear layer followed by the softmax function on top of the last layer hidden state corresponding to the [CLS] token. For PalmTree we use an LSTM over the instruction embedding tokens generated; to obtain a classifier we attach a linear layer with the softmax on top of the last hidden states of the RNN. We fine-tune each model for 20 epochs, selecting the epoch with the best classification accuracy on validation.

Results Results for both compiler and optimization classification are shown in Table 3. The compiler classification task is the only task where BinBert-FT is slightly worse than BinBert-MLM-FT, we believe that this is due to the fact that recognising a compiler signature is a syntactic task. PalmTree-FT has the worse performance among all fine-tuned models. Again we can see that pre-training is important as all the from-scratch models perform consistently worse than their fine-tuned counterpart. Results are similar for the optimization classification task.

6.2.3 Strand Recovery

Fine-tune Process and Dataset The dataset is made of 9267 basic blocks, each block contains at least 5 disjoint strands (i.e. the strands do not overlap on instructions). We split the dataset into 7412 training samples and 927 samples for both validation and test. We model strand recovery as instruction classification; given the instructions of a basic block and the final instruction of one of its strands, we aim at classifying the other instructions as either belonging to the same strand as the marked instruction or not. To mark an instruction we surround it with a special token.

We fine-tune BinBert and BinBert-MLM by attaching a classification head on top of the last layer hidden states of the first token of each instruction; the network will output 1 if an instruction is part of the strand to be recovered and 0 otherwise. As for previous tasks, we use an LSTM on top of PalmTree and we put a classification head on the hidden states of the first token of each instruction. We fine-tune each model for 20 epochs, selecting the epoch with the best classification accuracy on validation.

Results We reported precision, recall, and F1-score of the positive class in Table 3. We can see that the best performing model is BinBert-FT, it markedly surpasses PalmTree-FT on all the metrics considered. BinBert-MLM-FT is the second best model, it achieves the same precision as BinBert-FT but smaller recall. We believe the reason to be the execution-awareness of BinBert-FT that allows the model to recover more instructions in a strand.

6.2.4 Strand Execution

Fine-tune Process and Dataset We created a dataset for the strand execution task by taking strand-symbolic execution pairs, assigning concrete values to input variables, and evaluating its concrete output. In particular, we randomly assign values between 0 and 100 to input variables and we take only strands whose output is not greater than 200. Our dataset only contains strands computing the value of a register or a predicate of a conditional branch. The dataset is composed of 40000 training samples, and 5000 validation and test samples (total 50k). Each sample is made up of strand instructions followed by concrete assignments of input variables and the query output variables (only in the case of register outputs); the label for such a sample is the concrete value of the output variable. For instance, consider the strands `mov eax, dword ptr [rbp - 180]` `sub eax, 1` and the value 9 assigned to `dword ptr [rbp - 180]`. The corresponding sample will be `mov eax, dword ptr [rbp - 180]` `sub eax, 1` [SEP] `dword ptr [rbp - 180]` = 9 [SEP] `rax`. The network has to predict the value for register `rax`, in this case 8.

We model this problem as a sequence classification task, thus we used the same architectures used for the compiler provenance tasks (see Section 6.2.2).

Results We reported the accuracy obtained by all the models in Table 3. BinBert-FT achieves the best performances. PalmTree-FT performs poorly on this task, we believe that this is due to the fact that PalmTree is

a pure instruction embedding model so even when fine-tuned it cannot transfer to the upward neural architecture A sequence-related knowledge. On the contrary, both BinBert-FT and BinBert-MLM-FT are trained on sequences; however, BinBert-FT has the edge thanks to its execution-awareness. In this case, the pre-training has a great impact as we can see by the poor performance of BinBert-FS.

6.3 Extrinsic Tasks at Function Level

6.3.1 Similarity

Fine-tune Process and Dataset The function similarity dataset has exactly the same format as the similarity tasks at strand and block-level (see Section 6.2.1). In particular, the training dataset contains 18805 function pairs, while the validation set contains 1675 pairs. The test set is composed of 3627 functions, where 267 are queries. On average each query has 14 similar entities in the dataset.

On top of the embedding models, we use the same architectures that we employed to solve the other similarity tasks. As in [22] we truncate all functions to the first 150 instructions.

We compare our solution with Asm2vec [14] and Safe [22], which are solutions specifically crafted for the binary similarity problem. Asm2vec uses PV-DM [18] to simultaneously learn instruction and function embeddings, while Safe uses word2vec [24] to create instruction embeddings to be fed to a self-attentive neural network [21] for learning the final embeddings for functions.

Results Results are in Figure 7. BinBert achieves the best performances, beating also Asm2vec and Safe which are specific for the function similarity problem. PalmTree performs worse than BinBert, confirming our initial intuition about the drawbacks of the lack of context and the isolated instruction embeddings. Asm2vec is the worst performing model and is beaten by PalmTree (this aligns with the results in [19]). Safe has lower performances than PalmTree. This suggests that even if they are both based on fixed instructions embeddings, a transformer encoder is more capable of capturing an instruction semantic than word2vec. Also in this task, BinBert-FS and PalmTree-FS are below their corresponding fine-tuned versions (BinBert-FT and PalmTree-FT) thus highlighting the importance of pre-training. Finally, the impact of execution awareness can be seen in the advantage of BinBert-FT on BinBert-MLM-FT.

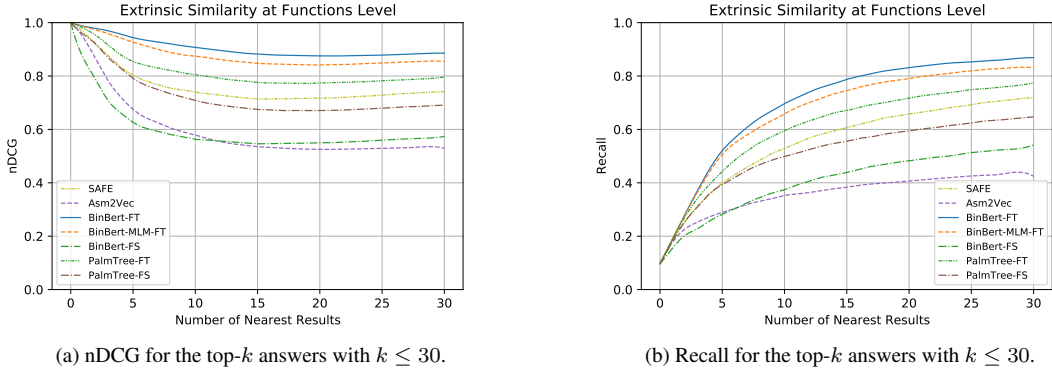


Figure 7: Results for the function similarity task.

Tasks			Models														
			BinBert-FT			BinBert-MLM-FT			BinBert-FS			PalmTree-FT			PalmTree-FS		
			Prec.	Rec.	nDCG	Prec.	Rec.	nDCG	Prec.	Rec.	nDCG	Prec.	Rec.	nDCG	Prec.	Rec.	nDCG
Similarity	Strands	top-10	60.9	63.1	79.9	53.3	55.9	72.9	31.3	32.4	49.2	46.4	48.0	65.3	38.4	39.8	56.8
		top-25	33.3	76.8	81.9	29.2	68.8	74.7	16.9	40.0	49.6	25.4	59.8	66.7	20.6	48.5	57.1
		top-40	23.3	81.1	83.0	20.7	73.8	76.2	12.2	44.1	50.7	18.1	64.6	68.2	14.8	52.9	68.6
	Basic Blocks	top-5	60.1	27.9	65.7	57.8	26.5	63.3	37.9	17.4	46.6	56.6	26.0	62.5	50.4	23.0	57.3
		top-10	44.6	39.9	55.5	42.5	37.9	53.4	24.6	22.0	36.0	39.6	35.4	51.1	33.7	30.0	54.5
		top-25	23.7	51.8	55.9	22.4	48.9	53.5	13.1	29.0	35.9	20.8	45.2	51.0	17.3	37.7	44.6
	Functions	top-5	94.3	45.3	95.7	92.4	44.2	94.3	57.2	25.8	66.2	83.7	39.5	88.0	77.5	36.0	83.0
		top-10	76.7	67.0	91.3	72.5	63.3	87.9	41.7	36.1	57.2	65.3	57.3	81.2	55.2	48.5	72.2
		top-25	43.2	85.1	87.8	41.2	81.5	84.7	24.2	50.6	55.7	36.7	74.2	78.0	29.8	61.8	67.7
			Accuracy			Accuracy			Accuracy			Accuracy			Accuracy		
Compiler Classification	Strands		78.4			78.9			59.3			73.5			70.6		
	Basic Blocks		75.4			76.1			62.8			71.3			67.3		
	Functions		89.1			89.4			72.9			85.6			81.3		
Optimization Classification	Strands		78.4			79.9			70.5			76.5			73.9		
	Basic Blocks		69.6			69.6			66.2			67.3			67.1		
	Functions		75.4			74.1			67.4			72.3			68.9		
Strand Execution			87.9			86.5			29.9			37.0			20.0		
			Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
Strand Recovery			96.9	98.7	97.8	96.9	98.4	97.6	65.7	29.9	41.1	81.1	75.5	79.2	80.1	56.7	66.3

Table 3: Extrinsic evaluation results.

6.3.2 Compiler Provenance

Fine-tune Process and Dataset The compiler provenance dataset has the same structure as the strand and block compiler provenance tasks. It contains 39017 functions in the training set and 4877 functions in the validation and test set.

We used the same architectures that we employed to solve the other compiler provenance tasks. We truncate all the functions to the first 150 instructions (as done in [22]).

Results We reported the accuracy for different models in Table 3. Results are similar to the compiler provenance task at strand and block level, thus the same considerations hold. In Appendix 9.2 we report the confusion matrices for this experiment.

Results from the experimental evaluation confirm that BinBert is the current state-of-the art for assembly code models. It shows improvement over PalmTree and deep learning solutions specifically tailored for a certain task. The execution awareness has a marked impact on semantic tasks. Interestingly, on some syntactic tasks, execution awareness does not increase the performance of the model. Our evaluation highlights the great impact of pre-training on downstream tasks with small size datasets.

7 Related Works

Binary code representation techniques can be subdivided into two main branches: manual features selection [20, 2, 26, 40, 28, 29] and unsupervised features extraction. Since our paper proposes an assembly model we discuss works that use or propose instruction embedding models. We first focus on the instruction embedding models proposed in the literature, categorising them according to the defining properties identified in the Background Section 2: the distributed representation learning used, the preprocessing of assembly instructions, and the

extraction methodology for assembly sequences. Finally, we discuss how these models have been used to solve binary analysis tasks.

7.1 Distributed Representation Learning

The distributed representation learning technique is the neural architecture and the training tasks used by the instruction models to create useful embedding vectors. The most common is word2vec [24]. Word2vec has been used, with minimal modifications, by Eklavya [9], SAFE [22] and others [42, 15, 23]. A notable difference is Asm2vec [14] which uses PV-DM [18], a variation of word2vec that simultaneously creates instructions and function embeddings. We remark that PV-DM cannot be fine-tuned. PalmTree [19], that we described in Section 2.1.3, uses a transformer.

7.2 Preprocessing of Assembly Instructions

The preprocessing is characterised by the substitution policy for information in the raw assembly instruction and the tokenization policy. We classify the substitution policies in *aggressive* or *light*. In an aggressive policy, lots of information contained in the assembly instructions are removed or changed. An example is DeepBinDiff [15] that replaces all constants and pointers with special tokens and renames registers according to their lengths in bytes (e.g. `ecx` becomes `reg4`). InnerEye [42], instead, replaces all constants, strings, and function names with special symbols. BinDeep [1] substitutes operands based on predefined categories (e.g. general register, direct memory reference, etc.). All the above policies are aggressive; indeed, InnerEye [42] wastes fundamental information that a library function call could bring and DeepBinDiff [15] loses register names that could be relevant, for instance, to understand the data flow.

Light preprocessing policies are applied by SAFE [22] and PalmTree [19] which keep small constant values and replace values above a predefined threshold with a special token. SAFE [22] replaces all memory addresses with the same token, PalmTree uses different tokens for generic memory locations and the ones pointing to strings.

Regarding tokenization, some works consider an entire instruction as a token [22] or split an instruction into opcode and operands [14], while others use a word-based tokenizer by splitting instructions on symbols (e.g. spaces or special characters) [30], [19]. The latter is a fine-grained tokenization strategy that is useful to reduce the vocabulary size and allow the upstream network to separately learn the semantics of each token (mnemonics, registers, etc.). No one used automatic tokenization strategies like WordPiece [39].

7.3 Binary Analysis Solutions using Embedding Models

Deep-learning-based solutions can be categorized into approaches using Graph Neural Networks (GNN), Recurrent Neural Networks (RNN), and Transformer architecture. [23, 11] use a GNN applied to function control flow graphs after transforming each block into a vector representation; this transformation is done by aggregating the instruction embeddings of each block. The resulting architecture is applied to the binary similarity [23], the compiler provenance problem [23], function naming problem [11]. Eklavya [9] and InnerEye [42] are examples of RNN-based solutions applied to the recovery of arguments used by a binary function [9], and basic-block similarity [42]. Another work is SAFE [22] which added a self-attention layer on top of an LSTM to solve the binary similarity problem. Among transformer-based solutions, we can find [3], [11] that apply transformers encoder-decoder architecture to recover function names from stripped binaries. Other works [41, 30] use transformer encoder for the binary similarity problem. In particular, [41] uses a transformer encoder to obtain basic blocks embeddings to be fed into a GNN which is trained in conjunction with a Convolutional Neural Network

(CNN) to produce a function embedding for binary similarity. Unfortunately, [41] does not release the code and the paper lacks of the details needed to reimplement their proposal. Trex [30] pre-trains a transformer encoder on function micro-traces and then uses this pre-trained model to produce embeddings for function similarity.

8 Conclusion

We presented BinBert, an execution-aware assembly language model. BinBert is trained on a big dataset of assembly strands and symbolic expressions. BinBert has shown state-of-the-art performance highlighting the relevance of execution-awareness. Our evaluation shows that BinBert is an encoder model that can be used to solve several tasks related to binary analysis using a fine-tune dataset of relatively small size. The generality of the model is an important strength and we believe that BinBert can be fruitfully applied to other downstream tasks as encoder layer of complex neural networks.

References

- [1] S. Alrabae, K.-K. R. Choo, M. Qbea'h, and M. Khasawneh, "Bindeep: Binary to source code matching using deep learning," in *Proceedings of the 2021 IEEE 20th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*, 2021, pp. 1100–1107.
- [2] D. Arp, M. Spreitzenbarth, M. Hubner, H. Gascon, and K. Rieck, "Drebin: Effective and explainable detection of android malware in your pocket," in *Proceedings of the 21st Annual Network and Distributed System Security Symposium (NDSS '14)*, 2014.
- [3] F. Artuso, G. A. Di Luna, L. Massarelli, and L. Querzoni, "Function naming in stripped binaries using neural networks," *CoRR*, vol. abs/1912.07946, 2019. [Online]. Available: <http://arxiv.org/abs/1912.07946>
- [4] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, "Enriching word vectors with subword information," *CoRR*, vol. abs/1607.04606, 2016.
- [5] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature verification using a "siamese" time delay neural network (nips'93)," in *Proceedings of the 6th Advances in Neural Information Processing Systems (NIPS '93)*, 1993, pp. 737–744.
- [6] J. Camacho-Collados and R. Navigli, "Find the word that does not belong: A framework for an intrinsic evaluation of word vector representations," in *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, 2016, pp. 43–50.
- [7] W. Chen, Y. Su, Y. Shen, Z. Chen, X. Yan, and W. Y. Wang, "How large a vocabulary does text classification need? a variational approach to vocabulary selection," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (ACL '19)*, vol. 1, 2019, pp. 3487–3497.
- [8] Y. Chen, Z. Shi, H. Li, W. Zhao, Y. Liu, and Y. Qiao, "Himalia: Recovering compiler optimization levels from binaries by deep learning," in *Proceedings of the 2018 Intelligent Systems and Applications (IntelliSys '18)*, K. Arai, S. Kapoor, and R. Bhatia, Eds., 2018, pp. 35–47.

- [9] Z. L. Chua, S. Shen, P. Saxena, and Z. Liang, “Neural nets can learn function type signatures from binaries,” in *Proceedings of the 26th USENIX Conference on Security Symposium (USENIX '17)*, 2017, pp. 99–116.
- [10] A. Conneau and G. Lample, “Cross-lingual language model pretraining,” in *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NeurIPS '19)*, 2019, pp. 7057–7067.
- [11] Y. David, U. Alon, and E. Yahav, “Neural reverse engineering of stripped binaries using augmented control flow graphs,” *Proceedings of the ACM on Programming Languages*, vol. 4, pp. 1–28, 2020.
- [12] Y. David, N. Partush, and E. Yahav, “Statistical similarity of binaries,” in *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '16)*, 2016, p. 266–280.
- [13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *CoRR*, vol. abs/1810.04805, 2018. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [14] S. H. H. Ding, B. C. M. Fung, and P. Charland, “Asm2vec: Boosting static representation robustness for binary clone search against code obfuscation and compiler optimization,” in *Proceedings of the 2019 IEEE Symposium on Security and Privacy (SP '19)*, 2019, pp. 472–489.
- [15] Y. Duan, X. Li, J. Wang, and H. Yin, “Deepbindiff: Learning program-wide code representations for binary diffing,” in *Proceedings of the 27th Annual Network and Distributed System Security Symposium (NDSS '20)*, 2020.
- [16] C. Fu, H. Chen, H. Liu, X. Chen, Y. Tian, F. Koushanfar, and J. Zhao, “Coda: An end-to-end neural program decompiler,” in *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NeurIPS '19)*, vol. 32, 2019, pp. 3703–3714.
- [17] Z. Hu, T. Chen, K.-W. Chang, and Y. Sun, “Few-shot representation learning for out-of-vocabulary words,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL '19)*, 2019, pp. 4102–4112.
- [18] Q. Le and T. Mikolov, “Distributed representations of sentences and documents,” in *Proceedings of the 31st International Conference on International Conference on Machine Learning (ICML'14)*, vol. 32, 2014, p. II–1188–II–1196.
- [19] X. Li, Y. Qu, and H. Yin, “Palmtree: Learning an assembly language model for instruction embedding,” in *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security (CCS '21)*, 2021, p. 3236–3251.
- [20] C. Liangboonprakong and O. Sornil, “Classification of malware families based on n-grams sequential pattern features,” in *Proceedings of the 2013 IEEE 8th Conference on Industrial Electronics and Applications (ICIEA '13)*, 2013, pp. 777–782.
- [21] Z. Lin, M. Feng, C. N. d. Santos, M. Yu, B. Xiang, B. Zhou, and Y. Bengio, “A structured self-attentive sentence embedding,” *CoRR*, vol. abs/1703.03130, 2017. [Online]. Available: <https://arxiv.org/abs/1703.03130>

- [22] L. Massarelli, G. A. D. Luna, F. Petroni, L. Querzoni, and R. Baldoni, “Function representations for binary similarity,” *IEEE Transactions on Dependable and Secure Computing*, pp. 1–1, 2021.
- [23] L. Massarelli, G. A. Di Luna, F. Petroni, L. Querzoni, and R. Baldoni, “Investigating graph embedding neural networks with unsupervised features extraction for binary analysis,” in *Proceedings of the 2019 Workshop on Binary Analysis Research (BAR ’19)*, 2019.
- [24] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS ’13)*, vol. 26, 2013, pp. 3111–3119.
- [25] Y. Otsubo, A. Otsuka, M. Mimura, T. Sakaki, and H. Ukegawa, “o-glassesx: Compiler provenance recovery with attention mechanism from a short code fragment,” *Proceedings 2020 Workshop on Binary Analysis Research (BAR ’20)*, 2020.
- [26] B. M. Padmanabhuni and H. B. K. Tan, “Buffer overflow vulnerability prediction from x86 executables using static analysis and machine learning,” in *Proceedings of the 2015 IEEE 39th Annual Computer Software and Applications Conference (COMPSAC ’15)*, vol. 2, 2015, pp. 450–459.
- [27] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Proceedings of the Advances in Neural Information Processing Systems (NEURIPS ’19)*, 2019, pp. 8024–8035.
- [28] J. Patrick-Evans, L. Cavallaro, and J. Kinder, “Probabilistic naming of functions in stripped binaries,” in *Proceedings of the 36th Annual Computer Security Applications Conference (ACSAC ’20)*, 2020, pp. 373—385.
- [29] J. Patrick-Evans, M. Dannehl, and J. Kinder, “XFL: extreme function labeling,” *CoRR*, vol. abs/2107.13404, 2021. [Online]. Available: <https://arxiv.org/abs/2107.13404>
- [30] K. Pei, Z. Xuan, J. Yang, S. Jana, and B. Ray, “Trex: Learning execution semantics from micro-traces for binary similarity,” *CoRR*, vol. abs/2012.08680, 2020. [Online]. Available: <https://arxiv.org/abs/2012.08680>
- [31] J. Pennington, R. Socher, and C. Manning, “GloVe: Global vectors for word representation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP ’14)*, 2014, pp. 1532–1543.
- [32] D. Pizzolotto and K. Inoue, “Identifying compiler and optimization options from binary code using deep learning approaches,” in *Proceedings of the 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME ’20)*, 2020, pp. 232–242.
- [33] E. C. R. Shin, D. Song, and R. Moazzezi, “Recognizing functions in binaries with neural networks,” in *Proceedings of the 24th USENIX Conference on Security Symposium (USENIX ’15)*, 2015, pp. 611–626.
- [34] Y. Shoshitaishvili, R. Wang, C. Salls, N. Stephens, M. Polino, A. Dutcher, J. Grosen, S. Feng, C. Hauser, C. Kruegel, and G. Vigna, “Sok: (state of) the art of war: Offensive techniques in binary analysis,” in *Proceedings of the IEEE Symposium on Security and Privacy (S&P ’16)*, 2016, pp. 138–157.

- [35] L. van der Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research*, vol. 9, no. 86, pp. 2579–2605, 2008.
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems (NEURIPS ’17)*, 2017, pp. 6000–6010.
- [37] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, “GLUE: A multi-task benchmark and analysis platform for natural language understanding,” in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 2018, pp. 353–355.
- [38] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, “Transformers: State-of-the-art natural language processing,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP ’19)*, 2020, pp. 38–45.
- [39] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. R. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. S. Corrado, M. Hughes, and J. Dean, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *CoRR*, vol. abs/1609.08144, 2016. [Online]. Available: <http://arxiv.org/abs/1609.08144>
- [40] X. Xu, C. Liu, Q. Feng, H. Yin, L. Song, and D. Song, “Neural network-based graph embedding for cross-platform binary code similarity detection,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS ’17)*, 2017, pp. 363—376.
- [41] Z. Yu, R. Cao, Q. Tang, S. Nie, J. Huang, and S. Wu, “Order matters: Semantic-aware neural networks for binary code similarity detection,” in *Proceedings of the 24th AAAI Conference on Artificial Intelligence (AAAI ’20)*, vol. 34, 2020, pp. 1145–1152.
- [42] F. Zuo, X. Li, P. Young, L. Luo, Q. Zeng, and Z. Zhang, “Neural machine translation inspired binary code similarity comparison beyond function pairs,” in *Proceedings of the 26th Network and Distributed Systems Security Symposium (NDSS ’19)*, 2019.

9 Appendix

9.1 Opcode and Operand Classes

In Table 4 and 5 we reported the classes that we have used to categorise instructions for the opcode and operand outlier detection tasks respectively. We performed a qualitative analysis based on the visualization of the clusters of opcodes learned by BinBert. To do so, we first used BinBert to convert each opcode into a vector, and we then applied t-SNE [35] to visualize opcodes in a two-dimensional space. Results are shown in Figure 8. Opcodes are well clustered according to their semantics. Specifically, we can identify two symmetric regions in which jumps, conditional move, and conditional set are split based on the condition checked (either positive or negative). We can also identify a region containing operations performing multiplications (`imul`, `leaq`) and divisions (`divq`), and other arithmetic operations (`add`, `sub`). Another interesting example is given by the region containing `ret`, `pop`, `push` and `call`: they all manipulate the stack.

Type	Opcodes
Data Movement	mov, push, pop, cwtl, cltq, cqto, cqtd
Unary Operations	inc, dec, neg, not
Binary Operations	add, sub, imul, xor, or, and, lea, leaq
Shift Operations	sal, sar, shr, shl
Special Arithmetic Operations	imulq, mulq, idivq, divq
Comparison and Test Instructions	cmp, test
Conditional Set Instructions	sete, setz, setne, setnz, sets, setns, setg, setnle, setge, setnl, setl, setnge, setle, setng, seta, setnbe, setae, setnb, setbe, setna
Jump Instructions	jmp, je, jz, jne, jnz, js, jns, jg, jnle, jge, jnl, jl, jnge, jle, jng, ja, jnbe, jae, jnb, jb, jnae, jbe, jna
Conditional Move Instructions	cmovz, cmovne, cmovnz, cmovs, cmovns, cmovg, cmovnl, cmovge, cmovle, cmovng, cmova, cmovnbe, cmovae, cmovnb, cmovb, cmovnae, cmovbe, cmovna
Procedure Call Instructions	call, leave, ret, retn
Floating Point Arithmetic	fabs, fadd, faddp, fchs, fdiv, fdivp, fdivr, fdivrp, fiadd, fidivr, fimul, fisub, fisubr, fmul, fmulp, fprem, fpreml, frndint, fscale, fsqrt, fsub, fsubp, fsubr, fsubrp, fextract

Table 4: Opcode classes used to categorise instructions for the opcode outlier detection task.

Type	Operand 1	Operand 2	Operand 3
none	-	-	-
cnst	immediate	-	-
reg	register	-	-
ref	memory	-	-
reg-reg	register	register	-
reg-cnst	register	immediate	-
reg-ref	register	memory	-
ref-reg	memory	register	-
ref-cnst	memory	immediate	-
tri	any	any	any

Table 5: Operand classes used to categorise instructions for the operand outlier detection task.

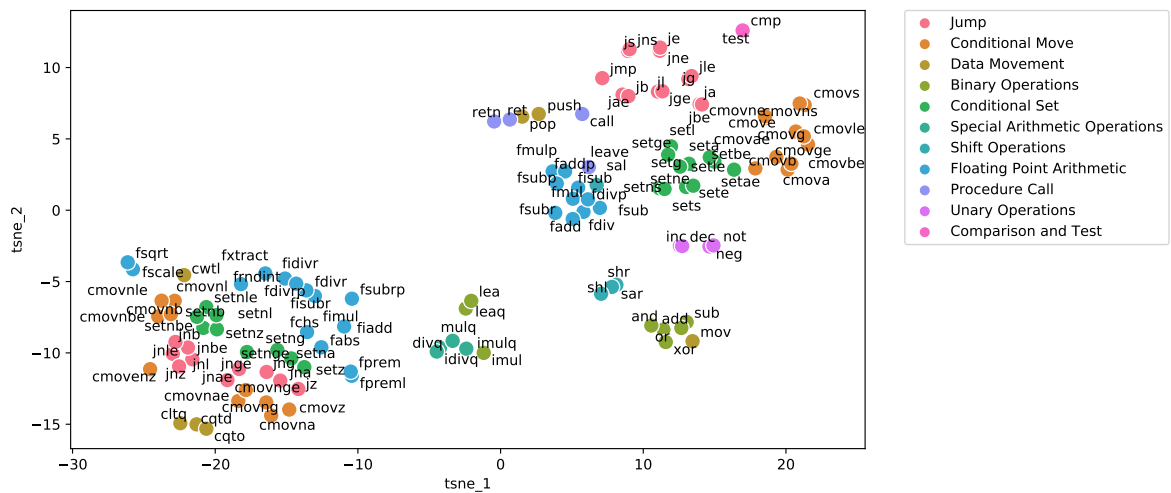
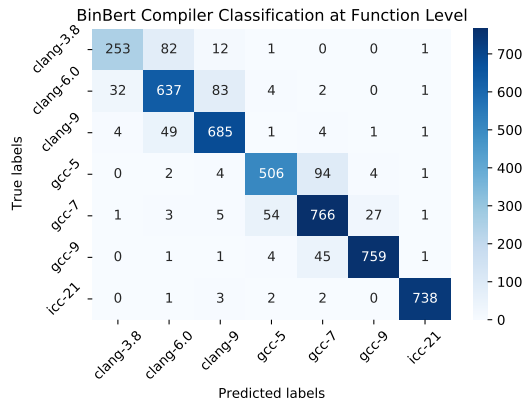
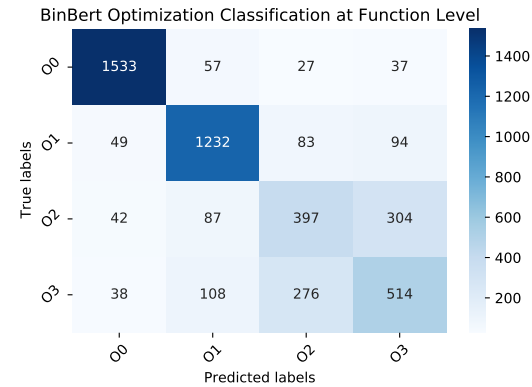


Figure 8: Visualization of opcode embeddings in a two-dimensional space with t-SNE.



(a) Confusion matrix for the compiler classification task at function level with BinBert.



(b) Confusion matrix for the optimization classification task at function level with BinBert.

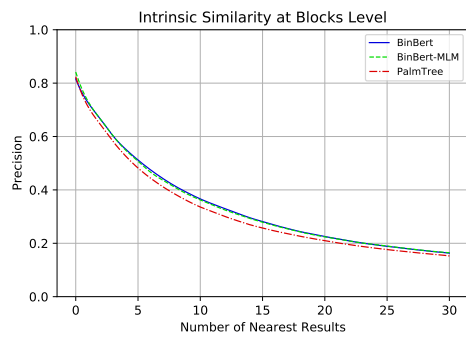
Figure 9: Results for the **extrinsic compiler provenance** task.

9.2 Figures for Compiler Provenance at Function Level

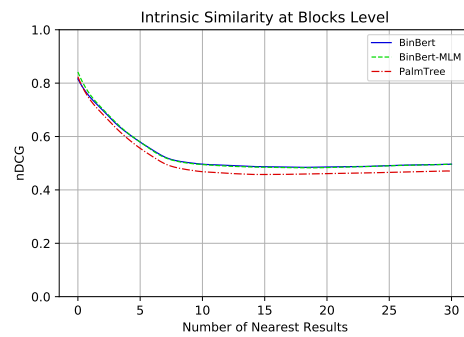
In Figure 9 we reported the confusion matrices obtained by using BinBert on the compiler and optimization classification tasks. We can observe that BinBert can clearly distinguish among different compiler families and it only gets confused with different versions in the same family. Similar behavior can be observed in the optimization classification task. It is easier for BinBert to distinguish optimized (O1, O2, O3) vs unoptimized code (O0) than to recognize the specific optimization level used to compile it.

9.3 Figures for Extrinsic and Intrinsic Block Similarity

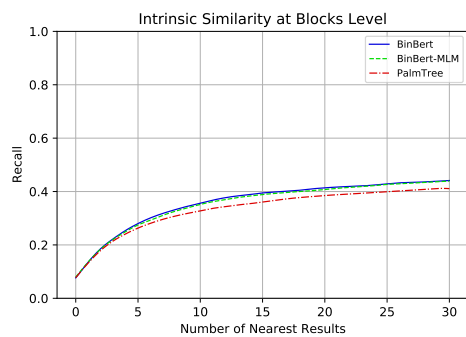
We reported the results of the intrinsic and extrinsic block similarity tasks in Figures 10 and 11. BinBert shows slightly higher performances with respect to PalmTree at the block level.



(a) Precision for the top- k answers with $k \leq 30$.

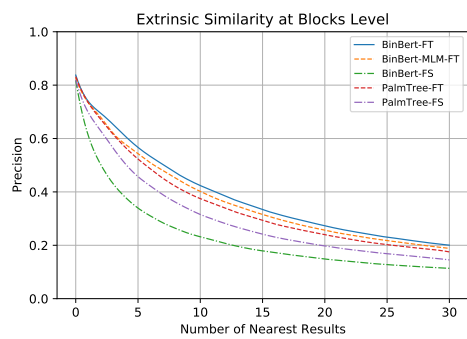


(b) nDCG for the top- k answers with $k \leq 30$.

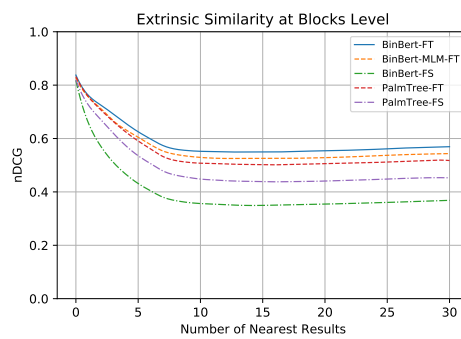


(c) Recall for the top- k answers with $k \leq 30$.

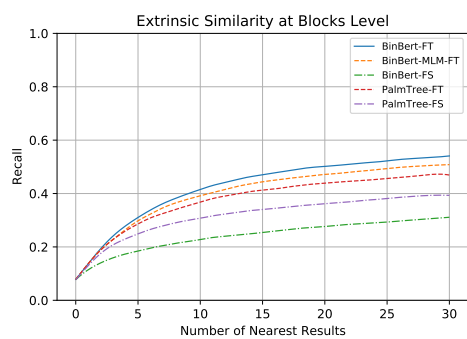
Figure 10: Results for the **intrinsic block similarity** task. Database of 6460 CFG basic blocks, average on 505 queries.



(a) Precision for the top- k answers with $k \leq 30$.



(b) nDCG for the top- k answers with $k \leq 30$.



(c) Recall for the top- k answers with $k \leq 30$.

Figure 11: Results for the **extrinsic block similarity** task. Database of 6460 CFG basic blocks, average on 505 queries.