

# A Transformer-based Function Symbol Name Inference Model from an Assembly Language for Binary Reversing

---

**HyunJin Kim**, JinYeong Bak, Kyunghyun Cho, Hyungjoon Koo

AsiaCCS '23

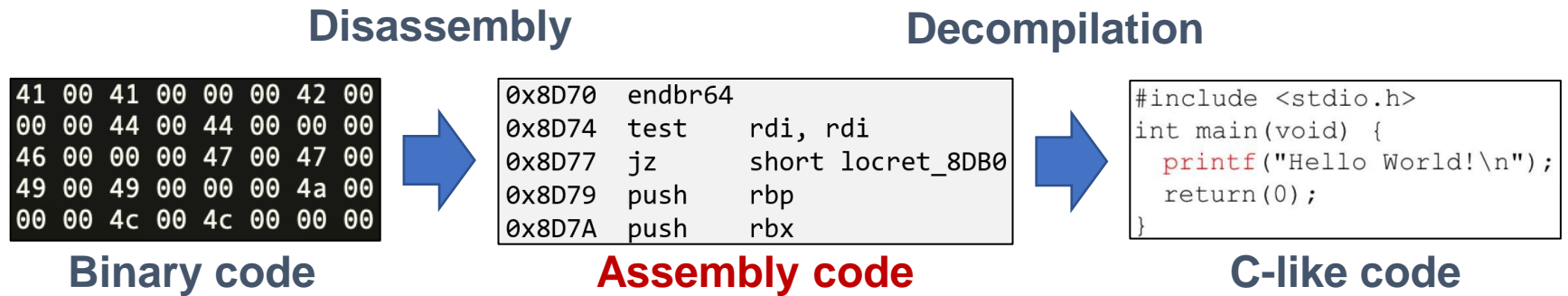
khyunjin1993@skku.edu



# Binary Reverse Engineering

---

- Gain the semantics from machine code



- Assembly instructions are made up of an opcode and operand(s)
- Opcode specifies an operation to perform
- Operand specifies data or memory location

# Missing Pieces

---

- Decompiler allows reversing engineers to infer code semantics by converting binary code to source in a C-syntax-like format
- Most high-level information is lost during the compilation
- It is infeasible to recover certain high-level information (e.g., variable name, **function name**, variable type, # of parameters)

```
1 void FUN_00108d70(void *param_1) {  
2     void *pvVar1;  
3     if (param_1 == (void *)0x0) {  
4         return;  
5     }
```

**Lost  
information**



# Binary Reversing

Difficult to comprehend a machine language (assembly)

→ A function often conveys a meaningful chunk with a name

## Assembly Language (Machine Language)

```
...  
0x8D70 endbr64  
0x8D74 test    rdi, rdi  
0x8D77 jz     short locret_8DB0  
0x8D79 push   rbp
```



A well-developed program maintains a well-described function label!

Binary

```
0x8DA5 jnz    short loc_8DB0  
0x8DA5 add    rsp, 8  
0x8DA9 pop    rbx  
0x8DAA pop    rbp  
0x8DAB retn
```



# Existing Work

Many researchers tried to analyze function names in a binary file

- Neural Reverse Engineering of Stripped Binaries using Augmented Control Flow Graphs (NERO) [David et al., OOPSLA 2020]

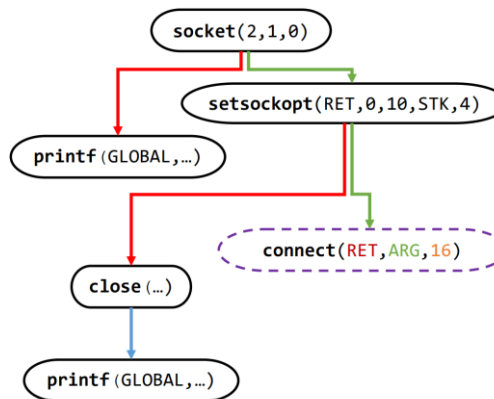
## ① Extract Assembly



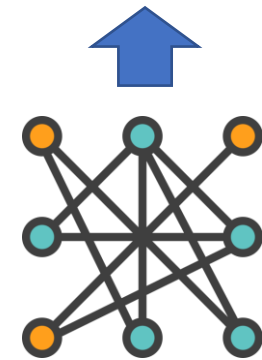
Binary

```
0x8D70 endbr64
0x8D74 test    rdi, rdi
0x8D77 call    socket
...
0x8DA5 call    setsockopt
...
0x8DAA call    connect
...
```

## ② Call Sites Graph



## Function name prediction

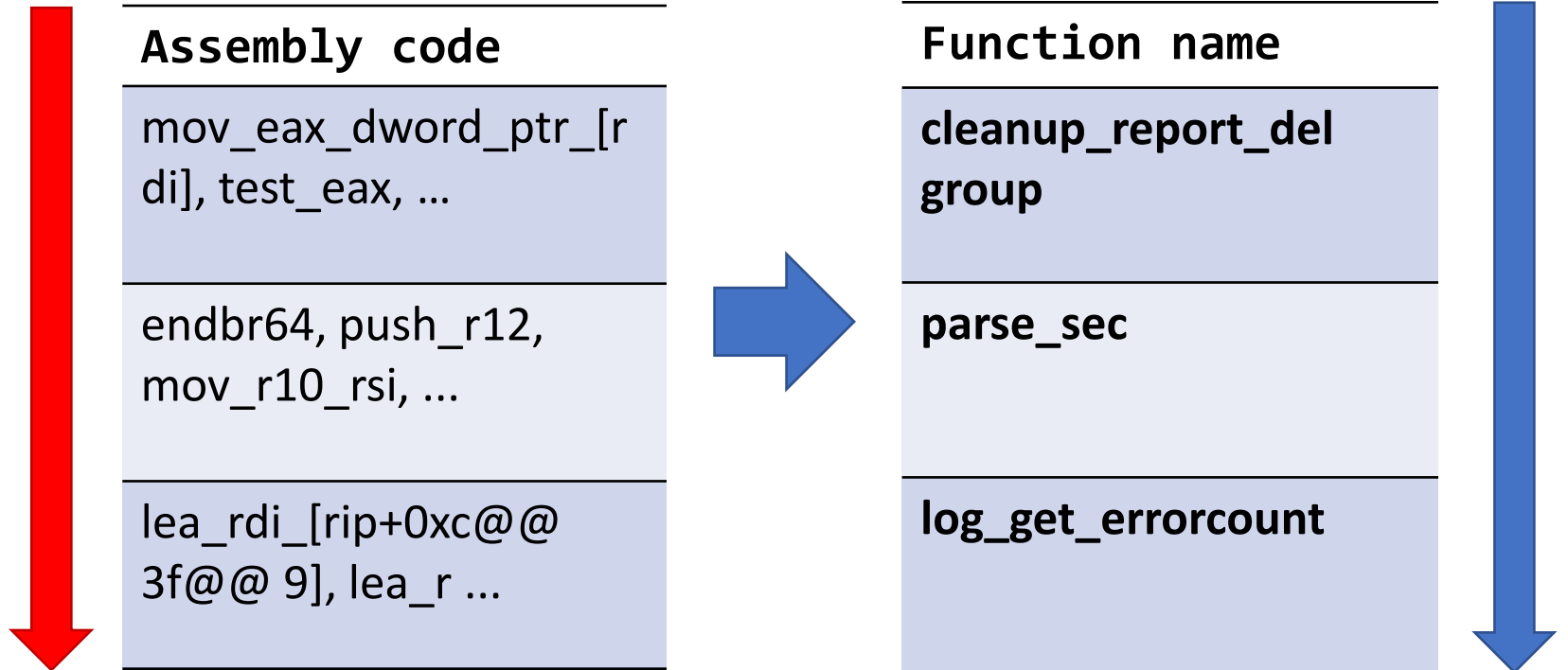


## ③ Graph Convolutional Network

# Goal

---

A series of function names allows reversers to **gain a quick overview** of a binary if they could be accurately inferred



# A Problem in NLP?

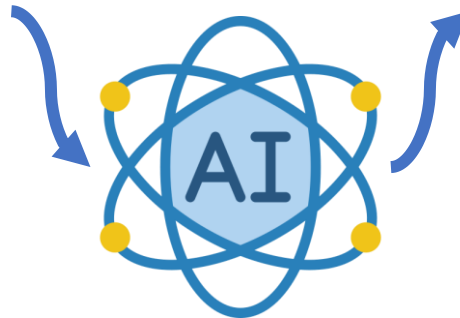
---

- Input and output are similar in terms of language
- The problem can be viewed as a **language translation task**

**Input  
(One Language)**

endbr64, test\_rdi\_rdi, jz  
short\_locret\_8DB0, ..., retn

안녕하세요  
(annyeonghaseyo)



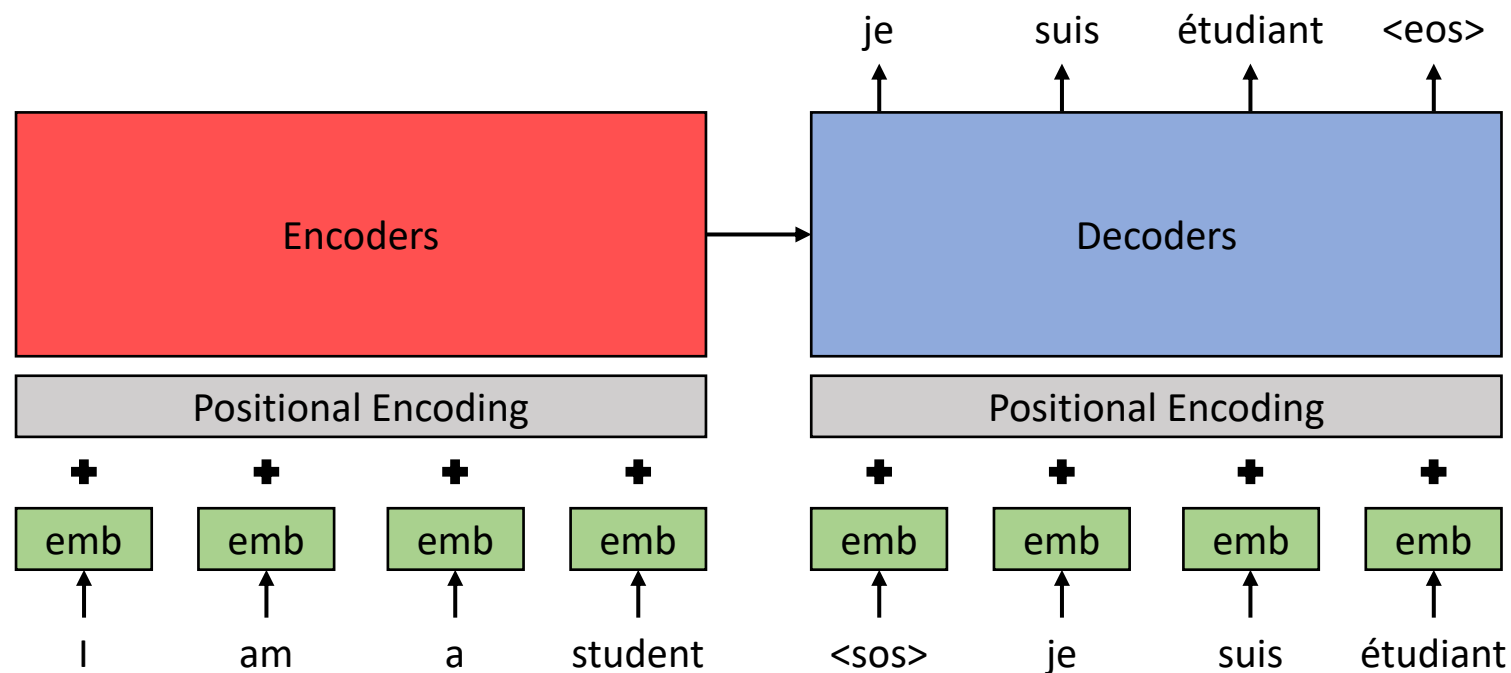
**Output  
(Another Language)**

ipc\_sem\_free\_info

Hello

# Transformer Language Model

- **Encoder** understands the meaning of input words (English)
- **Decoder** generates a sequence of words as output (French)





# Motivation

---

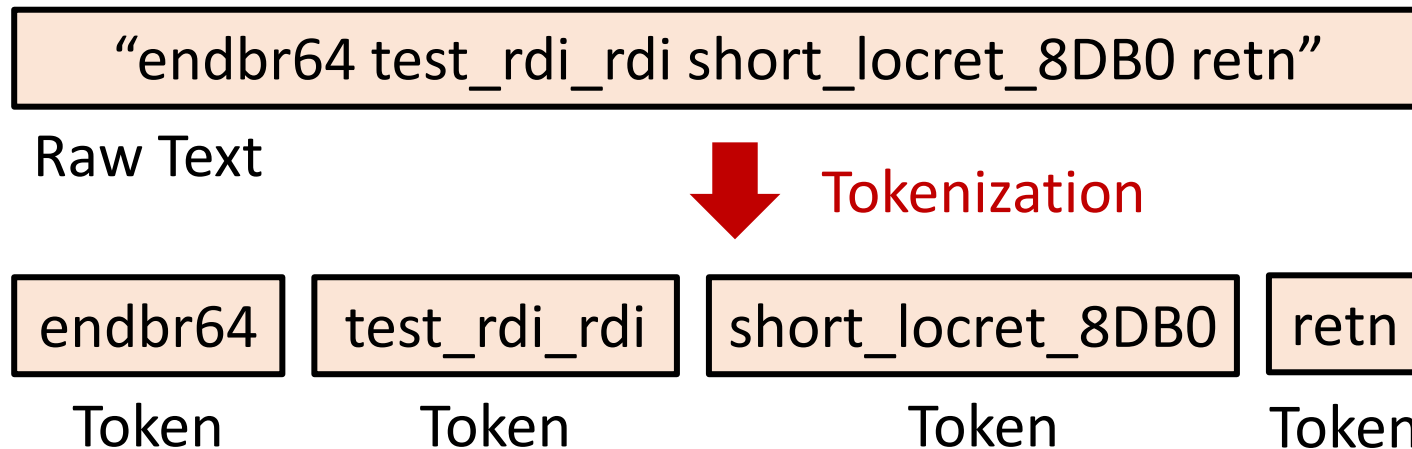
- Transformer: one of the breaking-through architectures in NLP
  - Widely adopted to translate between two languages
  - Naive Transformer structure struggles to learn the underlying semantics of machine instruction
  - Explore optimal ways of feeding data considering the structure of a language model
- **Design a model AsmDepictor tailored to a function name prediction task**

# Assembly Tokenization

---

Importance of a tokenization method

- Tokenization is splitting the raw text into small chunks of words or sentences, called tokens
- Tokenization determines how words are split, which affects the size of the vocabulary



# Assembly Tokenization (Instruction)

---

Challenges in instruction tokenization (e.g., `sub_rsp_0x50`)

- Numerous *sparse* vocabularies (e.g., Address, Immediate)
- 400 million tokens with 4 bytes of immediate value

All possible Opcode and Operand combinations

- Diverse outcomes in vocabulary combinations
- **Generates a tremendous of vocabularies (compared to NLP)**
- High computational cost with larger datasets

Tokenization Methods	Large Dataset (approx. 400,000 functions)	
	Vocabulary Size	Parameter Size
Instruction	3,253,394	1,012,857,350

# Assembly Tokenization (Byte Pair Encoding)

---

BPE tokenization for assembly instructions

Big Code!= Big Vocabulary: Open-Vocabulary Models for Source Code [ICSE, 2020]

- BPE calculates a sub-word frequency
- Widely adopted in the field of NLP (Natural Language Processing)

➤ **Significant reduction with a model parameter size**

sub\_rsp\_0x50 ➡ sub\_rsp\_0x5, 0

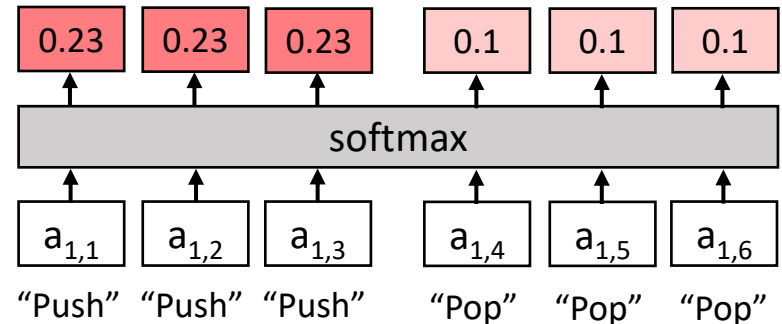
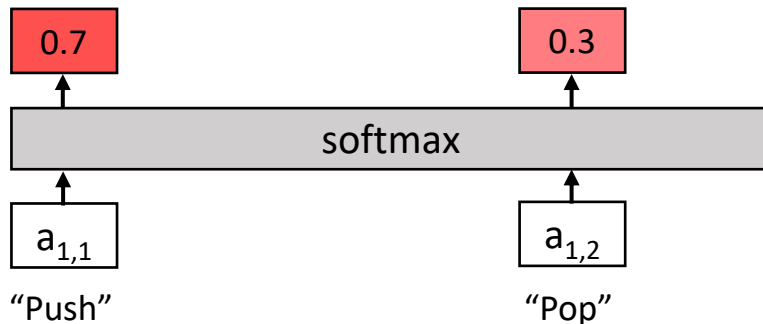
Tokenization Methods	Large Dataset (approx. 400,000 functions)	
	Vocabulary Size	Parameter Size
Instruction	3,253,394	1,012,857,350
<b>Byte Pair Encoding</b>	<b>9,923</b>	<b>40,004,102</b>

# Problem of Attention for Our Purpose

---

## Problem

- Frequent appearance of a machine instruction may convey important information (e.g., opcodes)
- Softmax reduces the collective output of duplicated words that might carry significant information

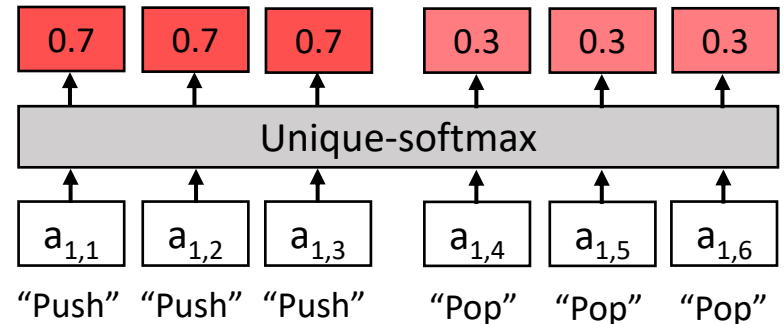
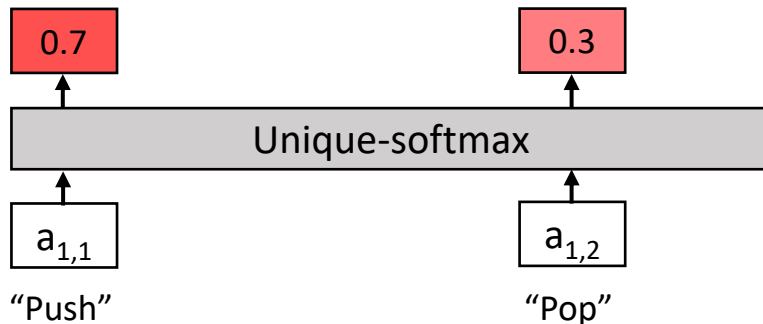


# Our Approach (1): Unique-softmax

---

## Unique-softmax

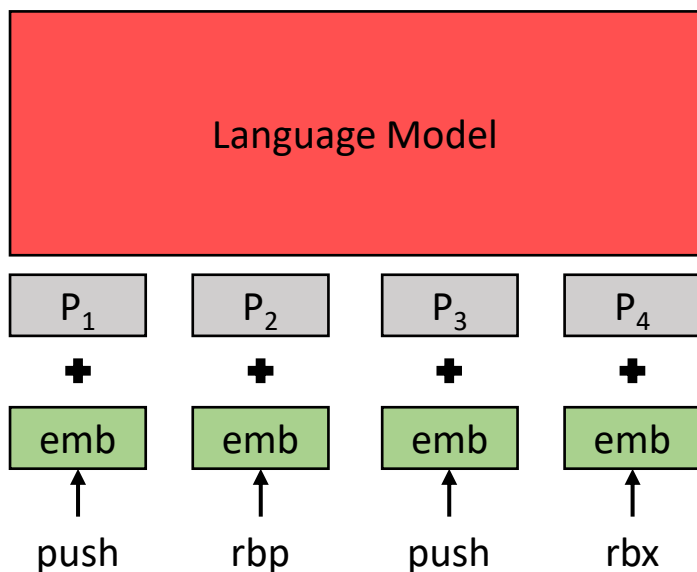
- Enables to calculate probability within a **unique set** of word
- Retain the collective output of duplicates  
(+ 3.17% on F1-score)



# Problem of Positional Encoding

## Problem

- Assigns a pre-defined value to give positional information
- They **cannot learn the representation of the position**
- Position information only added at the first layer



$$P(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

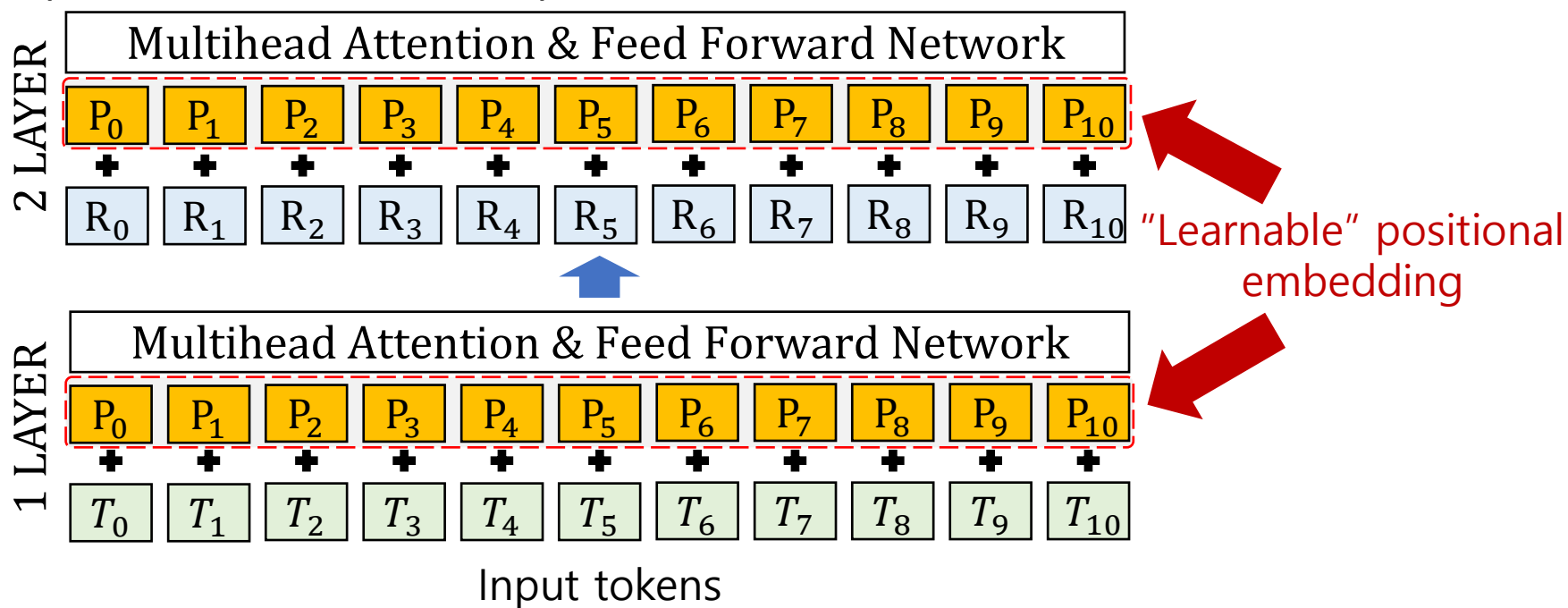
$$P(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

$i$  : dimension of the positional encoding vector  
 $pos$  : position of the token

# Our Approach (2): Per-layer Positional Embedding

## Per-Layer Positional Embedding

- Learns the positional representation
- Feed the order of machine instructions to every layer (+17.93% on F1-score)



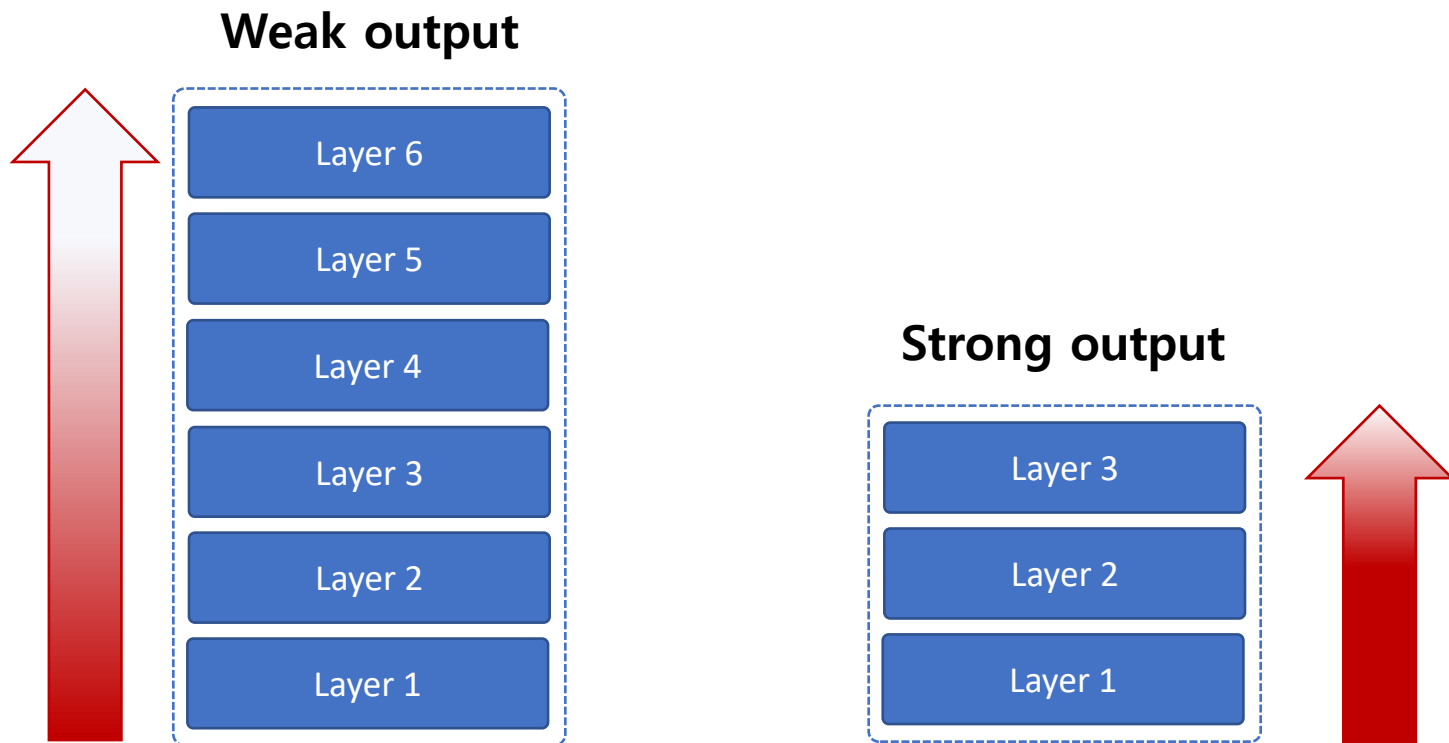


# Our Approach (3): Layer Reduction

---

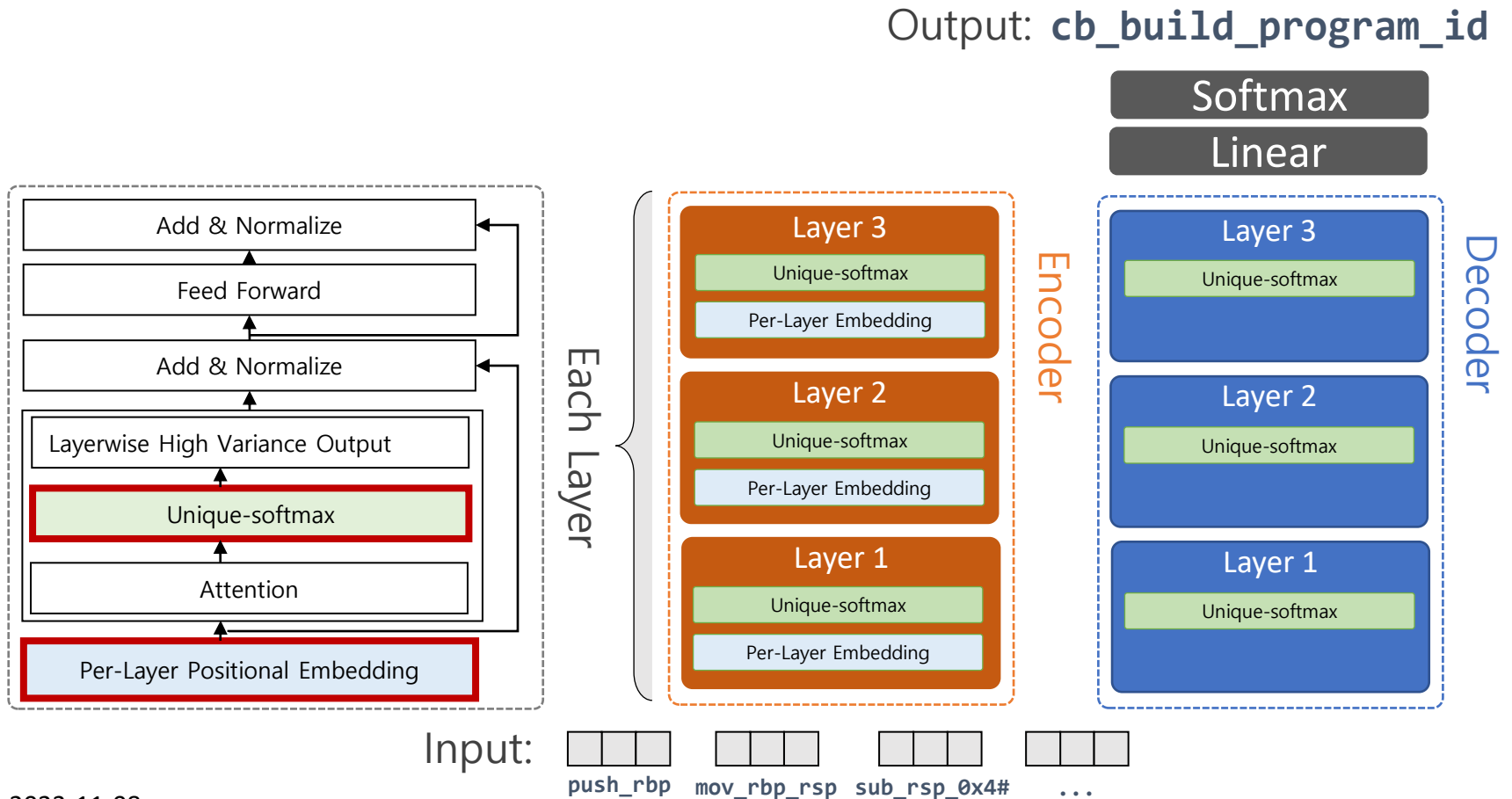
## Layer Reduction

- Reducing the layer prevents values diminishing at the upper layers (+5.64% on F1-score)



# AsmDepictor: Putting It All Together

## Unique-softmax + Per-Layer Positional Embedding + Layer reduction



# Experiment Settings

---

## Data

- Small dataset: 67,000 functions from NERO
- Large dataset: 400,000 functions from common packages from popular Ubuntu Linux distributions + Small dataset

## Metric

- F1-score: Token matching
- Rouge-l: Popular metric for longest common subsequences

## Model

- NERO: Neural reverse engineering of stripped binaries using augmented control flow graphs [David et al., OOPSLA 2020]
- Debin: Predicting Debug Information in Stripped Binaries [He et al., CCS 2018]

# Performance Comparison

- AsmDepictor surpasses existing SOTA models

		Model	Precision	Recall	F1	Rouge-I
Trained with Small Dataset	{	Debin	5.73	5.66	5.66	5.87
		NERO	12.35	12.36	12.35	14.07
		AsmDepictor	57.13	57.17	<b>57.14</b>	<b>58.68</b>

+44%

- Performance improvement with a large dataset

Model	Precision	Recall	F1	Rouge-I
AsmDepictor (Small)	57.13	57.17	<b>57.14</b>	<b>58.68</b>
AsmDepictor (Large)	71.52	71.53	<b>71.52</b>	<b>73.75</b>

+15%

# Demonstrative Examples

---

- Example of predictions
- **Red** indicates correctly predicted tokens

Debugging Symbol	AsmDepictor (Large Dataset)	AsmDepictor (Small Dataset)	NERO	Debin
cb_build program_id	<b>cb_build</b> <b>program_id</b>	<b>cb_build</b> <b>program</b>	options_menu	cint_remove
close_stdout	<b>close_stdout</b>	<b>close_stdout</b>	<b>close_stdout</b>	<b>close_stdout</b>
write_file	<b>write_file</b>	<b>write_file</b>	process	to_rgip
xcalloc	<b>xcalloc</b>	<b>xcalloc</b>	<b>xcalloc</b>	alloc_common

# Conclusion

---

## Wrap-up

- Our work introduces AsmDepictor, an effective prediction framework for a function name from machine instruction
- AsmDepictor surpasses existing state-of-the-art models, with F1 scores up to four times higher
- Future work: Integration function name inference model with a large language model such as LLaMA
- For more details of code & dataset, please visit our GitHub repository <https://github.com/agwaBom/AsmDepictor>

Thank you

---