# SymLM: Predicting Function Names in Stripped Binaries via Context-Sensitive Execution-Aware Code Embeddings

Xin Jin
The Ohio State University
jin.967@osu.edu

Kexin Pei
Columbia University
kpei@cs.columbia.edu

Jun Yeon Won
The Ohio State University
won.126@osu.edu

Zhiqiang Lin
The Ohio State University
zlin@cse.ohio-state.edu

## Overview

Function name prediction in stripped binaries is a very useful but extremely challenging task. For this,

- We design a **function symbol name prediction** and **binary language modeling** framework, SymLM.

- We propose a **novel neural architecture** to jointly learn function semantics preserved in **execution behavior** and **calling context**.

- We evaluate SymLM with **1.4M** binary functions and show that it outperforms the state-of-the-art works by **up to 35% in F1 score** with better generalizability and obfuscation resistance.

- We show SymLM's component effectiveness and **practical use cases** with IoT firmware images.

## Background and Motivation

1. Commercial software (e.g., IoT firmware, browsers, and pdf readers) is usually **closed-source** and shipped in **stripped** binaries, whose **semantic information** (e.g., function names) is **missing**. Predicting function names helps reverse engineers **understand code semantics, identify malware/vulnerabilities**, etc.

2. Predicting function names is **very challenging**, because:

a. Semantic similar code can appear differently, e.g.,



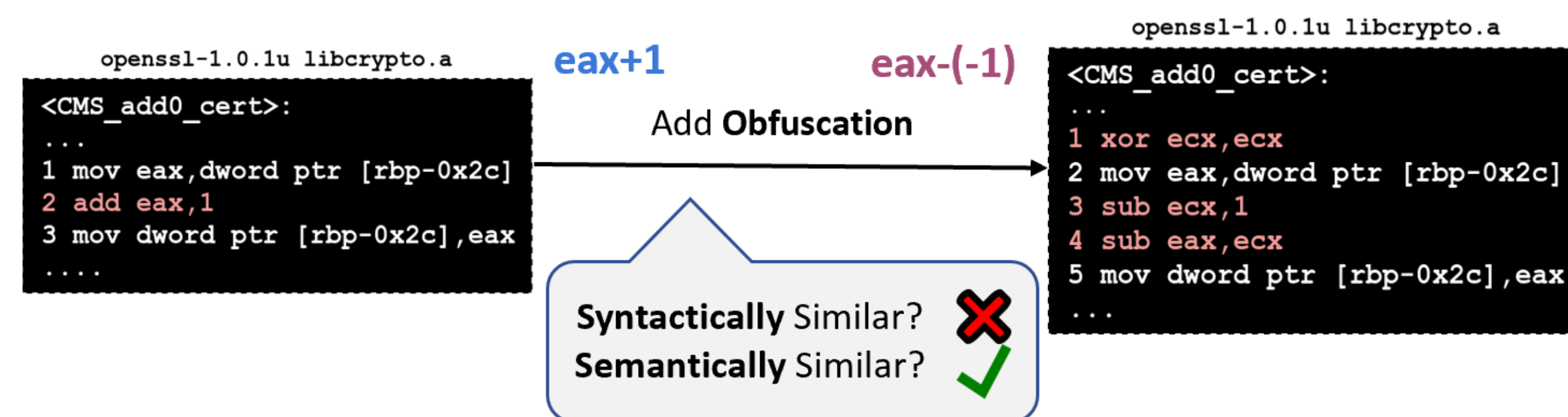Fig. 1: Semantically Similar but Syntactically Different Code

b. Function names are noisy, e.g., the probability that 2 developers give the same name for a function is 6.9%.

c. Modeling calling context is necessary, e.g.,

```
1  YY_BUFFER_STATE yy_scan_string (char *yystr) {
2      size_t _yybytes_len;
3      YY_BUFFER_STATE pyVar1;
4      _yybytes_len = strlen(yystr);
5      pyVar1 = yy_scan_bytes(yystr,_yybytes_len);
6      return pyVar1;
7  }
```

Fig. 2: Partial Semantics Preserved in Its Callees

**Our key observations**: predicting function names requires (1) learning semantics from execution behavior, (2) resolving NLP issues, and (3) modeling calling context.
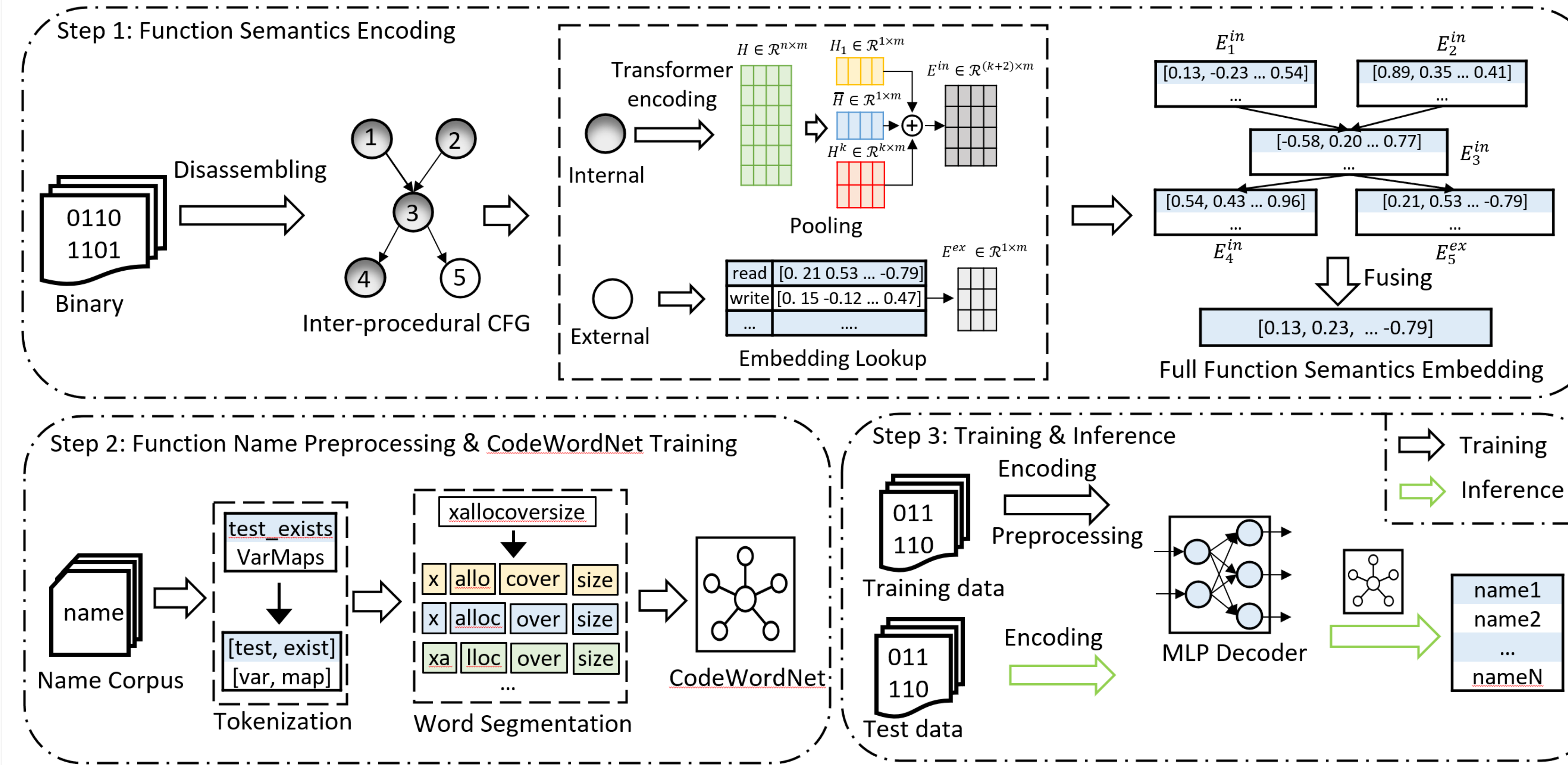
## Methodology



Fig. 3: System Workflow with Three Major Steps

**Step 1:** SymLM generates embeddings by fusing semantics of calling context and function instructions. It encodes internal functions by a pretrained model and external functions by an embedding lookup table.

**Step 2:** SymLM resolves NLP issues by tokenizing names into words, segmenting words by a unigram language model, and embedding words with CodeWordNet (consisting of 3 word embedding models).
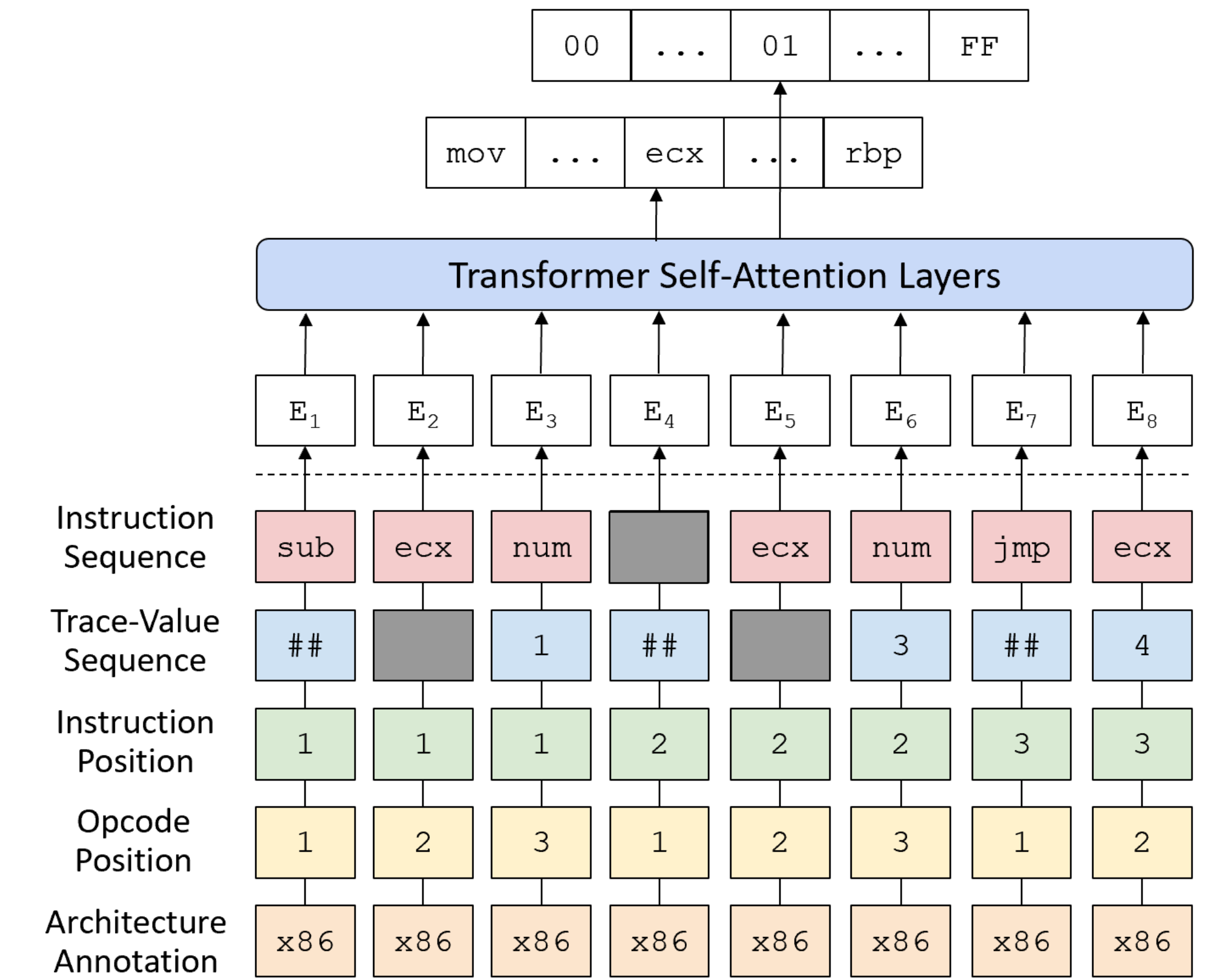
**Step 3:** SymLM updates weights of the pretrained model, embedding lookup layer, and MLP decoder in training, and predicts names of stripped binaries that semantically match ground truth with CodeWordNet.



Fig. 4: The Pretrained Model Used for **Transformer Encoding**. The pretraining tasks force the model to learn **execution behavior**.

## Evaluation and Results

**1. Evaluation Setup**: we built our datasets with 27 open-source projects, compiled into **16K binaries and 1.4 M functions** in 4 architectures, 4 optimizations, and 4 obfuscation options.

**2. Overall Performance**: 0.634 precision, 0.677 recall, and 0.655 F1 score on average.

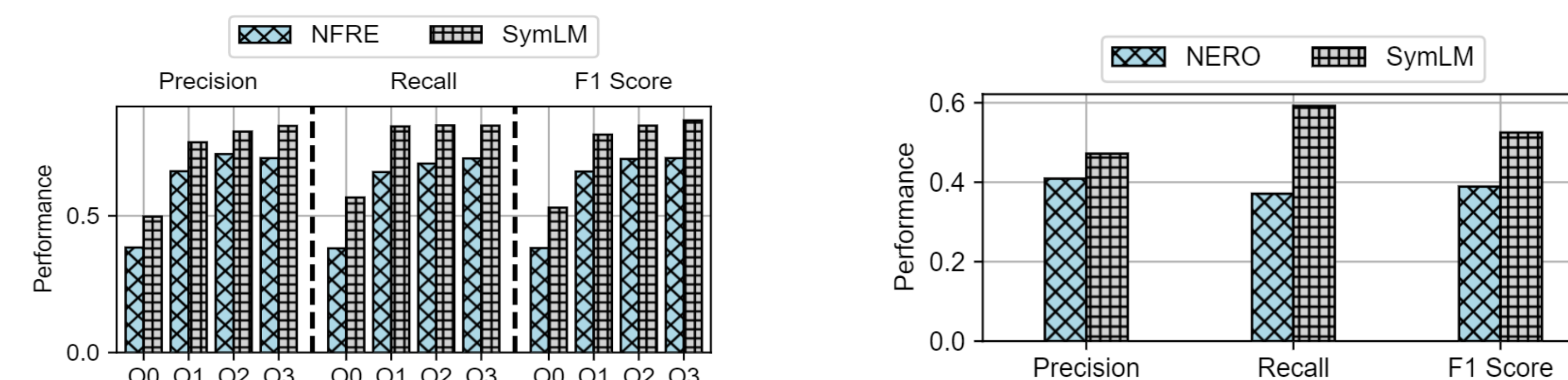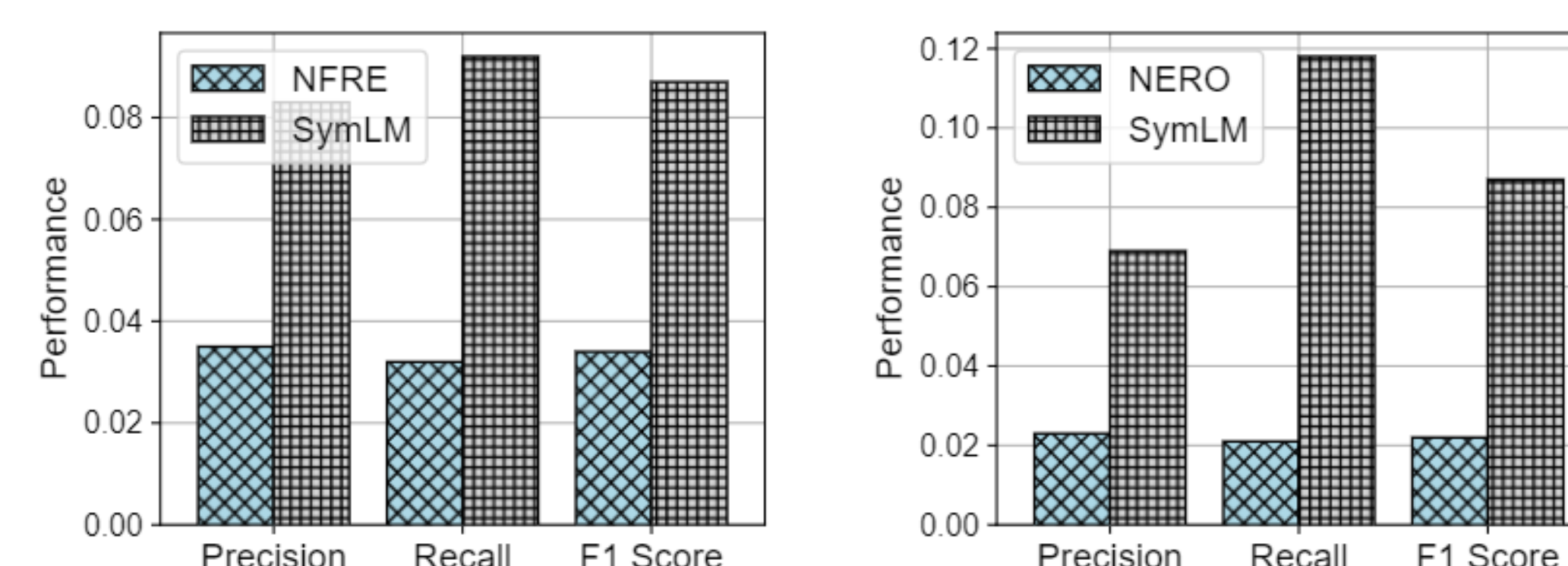**3. Comparison to the state-of-the-art works:**



Fig. 5: Baseline Comparison

**4. Generalizability Evaluation:**



(a) SymLM v.s. NFRE    (b) SymLM v.s. NERO

Fig. 6: Generalizability Comparison

**5. Obfuscation Resistance:**

Tab. 1: Performance (F1 Score) on Obfuscated (OBF) Binaries

| OBF | SymLM | NFRE |
|-----|-------|------|
| None | 0.806 | 0.595 |
| bcf | 0.757 (-6.1%) | 0.491 (-17.5%) |
| cff | 0.768 (-4.7%) | 0.445 (-25.2%) |
| sub | 0.726 (-9.9%) | 0.505 (-15.2%) |
| split | 0.788 (-2.2%) | 0.496 (-16.6%) |

## Use Case Study

We show the practical use case of SymLM with **8 32-bit ARM IoT firmware images** with 18% unseen IoT-specific name words (e.g., analog).
**Result**: 172/1062 names correctly predicted.

```
1   uint32_t read(uint32_t ulPin){
2       ...
3       if (pin == NC) uVar3 = 0;
4       else {
5           uVar2 = read_value(pin);
6           uVar3 = (uint32_t)uVar2;
7           if (uVar4 != 0xc) {
8               if ((uint) uVar4 < 0xc)
9                   return (uint)(uVar2 >> (0xcU - uVar4 & 0xff));
10                  return uVar3 << (uVar4 - 0xcU & 0xff));
11          }
12      }
13      return uVar3;
14  }
```

Fig. 7 Example Prediction. The ground truth names are *analogRead* and *abc_read_value*.

Scan the QR code or visit https://github.com/OSUSecLab/SymLM to find our code.