Motivation & Research Objective
oooo

Methodology
o

Results
o

Limitations
o

Conclusion & Future Work
o

# Comparing Natural Language Embeddings for Libc Functions as Rich Labels

## Bachelor defense

Ruben Triwari

Ludwig Maximilian University Munich

19, February 2025

Motivation & Research Objective
○○○○

Methodology
○

Results
○

Limitations
○

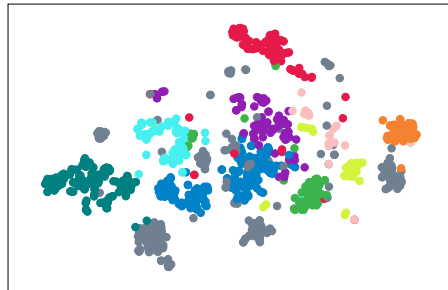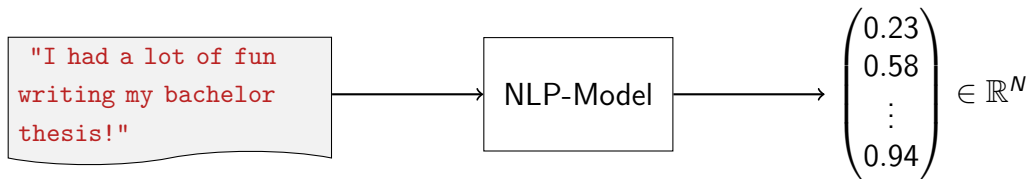Conclusion & Future Work
○

# Outline

Motivation & Research Objective

Methodology

Results

Limitations
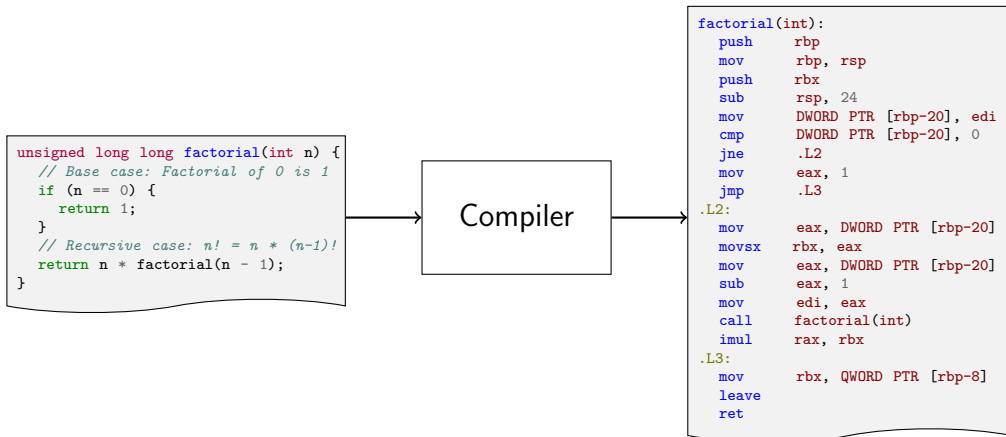
Conclusion & Future Work
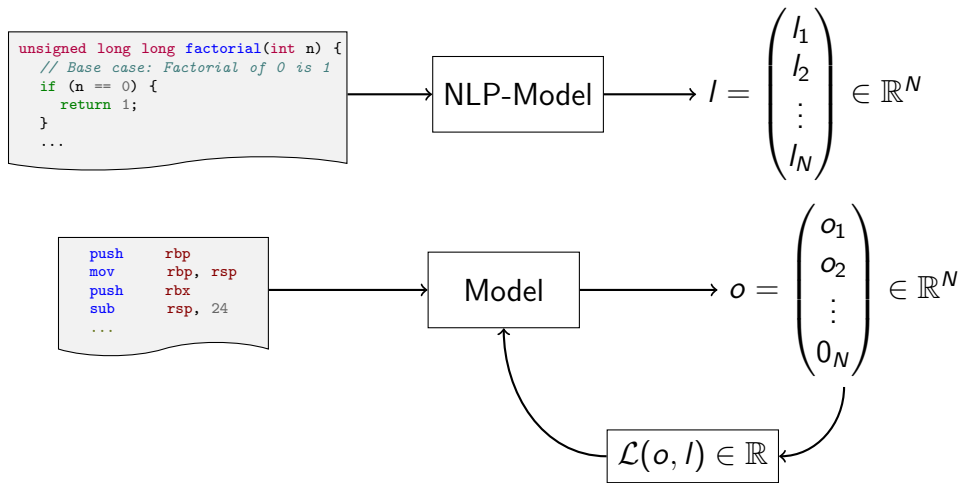
## Motivation



$\rightsquigarrow$ Encoding natural language was a huge factor in recent nlp advancements
$\rightsquigarrow$ Information described as a vector can be used in many downstream task
$\rightsquigarrow$ That motivates encoding binary code and describing them as a vector
$\rightsquigarrow$ That motivates using NLP tools to encode binary code

Motivation & Research Objective
○●○○
Methodology
○
Results
○
Limitations
○
Conclusion & Future Work
○

# Motivation

```
unsigned long long factorial(int n) {
    // Base case: Factorial of 0 is 1
    if (n == 0) {
        return 1;
    }
    // Recursive case: n! = n * (n-1)!
    return n * factorial(n - 1);
}
```

Compiler

```
factorial(int):
    push    rbp
    mov     rbp, rsp
    push    rbx
    sub     rsp, 24
    mov     DWORD PTR [rbp-20], edi
    cmp     DWORD PTR [rbp-20], 0
    jne     .L2
    mov     eax, 1
    jmp     .L3
.L2:
    mov     eax, DWORD PTR [rbp-20]
    movsx   rbx, eax
    mov     eax, DWORD PTR [rbp-20]
    sub     eax, 1
    mov     edi, eax
    call    factorial(int)
    imul    rax, rbx
.L3:
    mov     rbx, QWORD PTR [rbp-8]
    leave
    ret
```

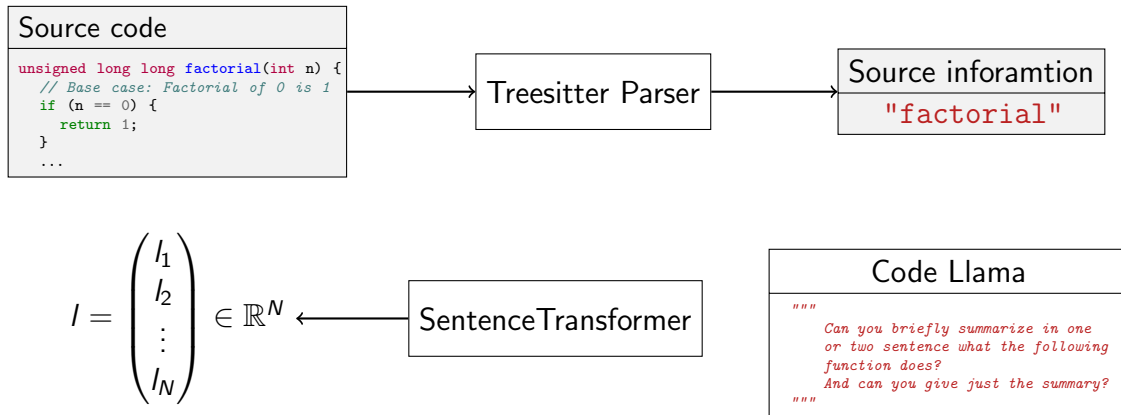⤳ Compiler removes important inforamtion in natural language

# Motivation

## Research Objectives

▶ Compare diffrent approaches generating an Embedding with NLP tools
  1. Embed function names with SentenceTransformer
  2. Embed function comments with SentenceTransformer
  3. Embed Code-Llama code summaries with SentenceTransformer
▶ Compare NLP approach to the existig Code2Vec Model
▶ Propose a new way comparing embedding spaces

Motivation & Research Objective
oooo

Methodology
•

Results
o

Limitations
o

Conclusion & Future Work
o

# Architecture



```
Source code
unsigned long long factorial(int n) {
    // Base case: Factorial of 0 is 1
    if (n == 0) {
        return 1;
    }
    ...
```

→ Treesitter Parser →

```
Source inforamtion
"factorial"
```

$$I = \begin{pmatrix} l_1 \\ l_2 \\ \vdots \\ l_N \end{pmatrix} \in \mathbb{R}^N$$

← SentenceTransformer ←

```
Code Llama
"""
    Can you briefly summarize in one
    or two sentence what the following
    function does?
    And can you give just the summary?
"""
```

# Denotational Semantics

HALLO

Motivation & Research Objective
oooo

Methodology
o

Results
o

Limitations
o

Conclusion & Future Work
●

HALLO

# Discussion