





# Comparing Natural Language Embeddings for Libc Functions as Rich Labels

Bachelor's Thesis Defense

Ruben Triwari

Ludwig Maximilian University Munich

February 18, 2025

# Outline

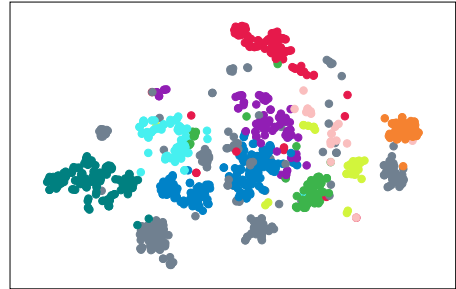
Motivation & Research Objective

Methodology

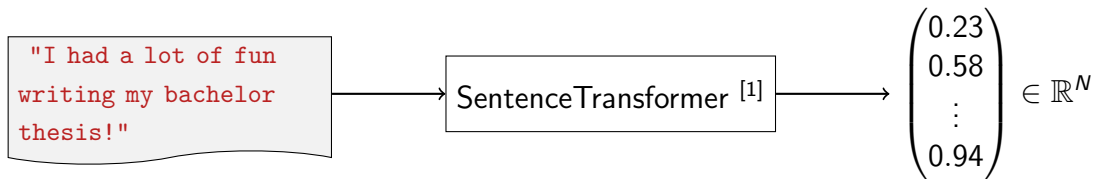
Results

Limitations

Conclusion & Future Work



# Motivation

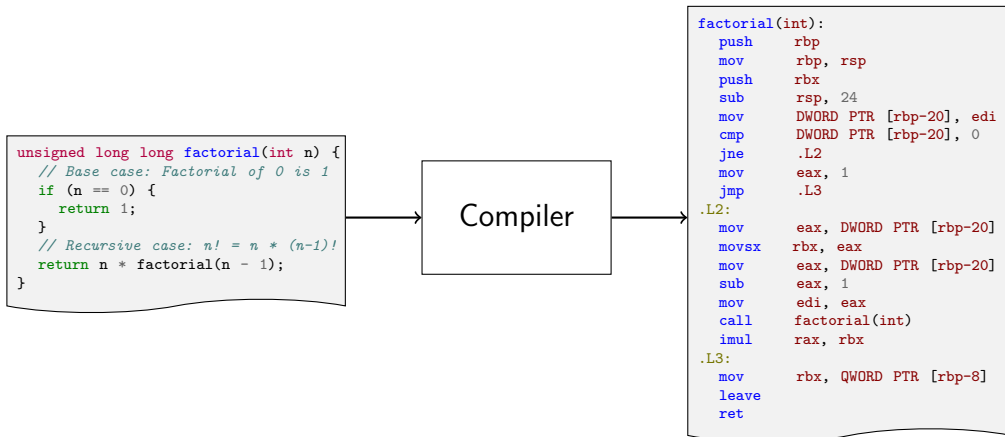


- ↪ Encoding natural language had an important role in recent NLP advancements
- ↪ Information described as a vector can be used in many downstream tasks
- ↪ That serves as an motivation for encoding binary code as vector
- ↪ That motivates using NLP tools to encode binary code

---

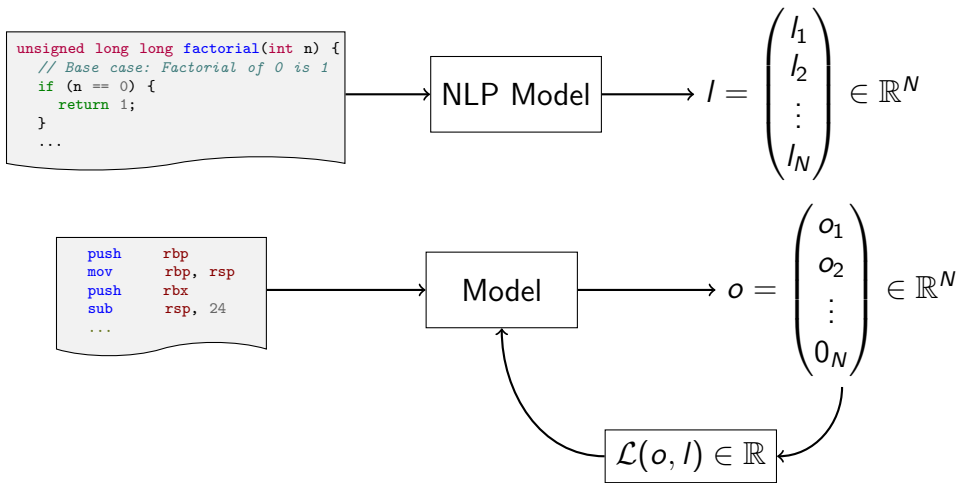
[1]Reimers and Gurevych: Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, EMNLP'19

# Motivation



⇒ Compiler removes all information that is in natural language

# Motivation



# Research Objectives

- ▶ Compare different approaches encoding additional information in the source code into machine readable format
  1. Embed function names with SentenceTransformer
  2. Embed function comments with SentenceTransformer
  3. Embed Code Llama <sup>[2]</sup> code summaries with SentenceTransformer

↪ Intuition is that Code Llama explanation will yield "good" embeddings
- ▶ Compare NLP approach to the existing Code2Vec <sup>[3]</sup> Model
- ▶ Propose a new way comparing embedding spaces.

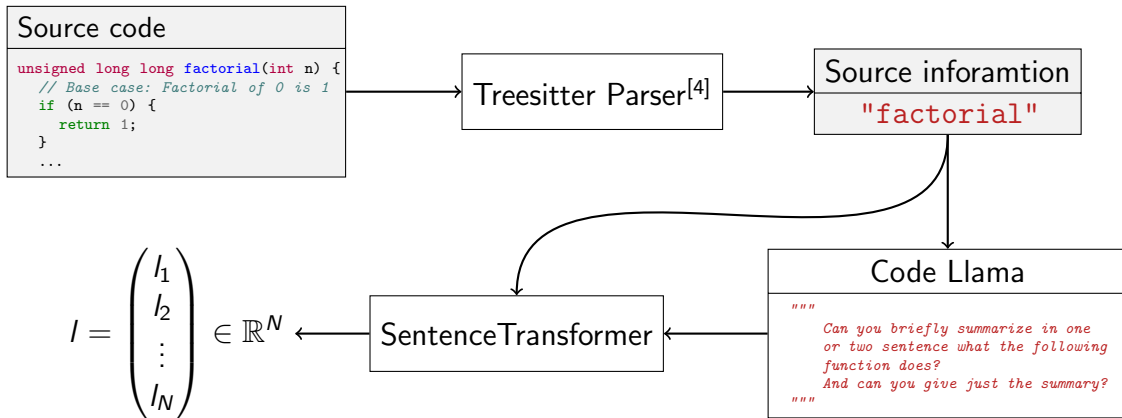
↪ To prove intuition

---

[2]Rozière et al.: Code Llama: Open Foundation Models for Code, 24

[3]Alon et al.: code2vec: Learning Distributed Representations of Code, POPL'19

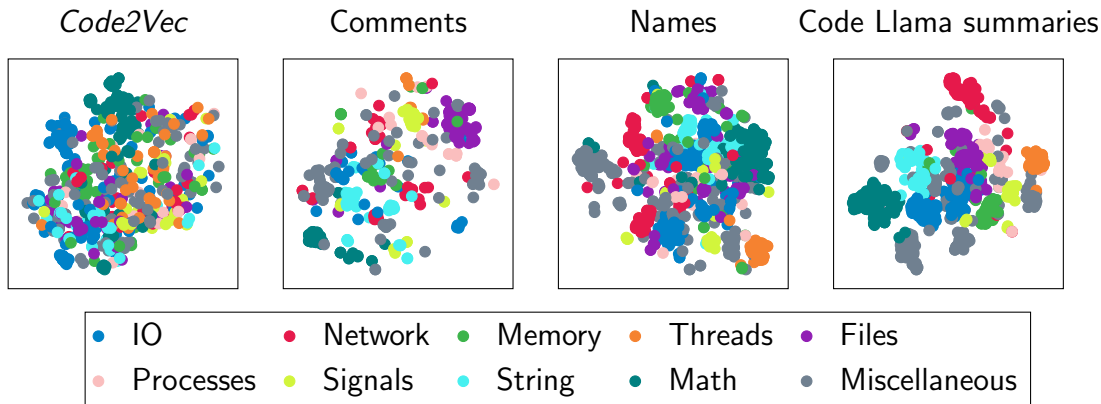
# Architecture



[4] Official website: <https://tree-sitter.github.io/tree-sitter/>



# Evaluation with t-SNE



**Figure:** Depicted are the *t-SNE* output vectors with perplexity  $P = 30$ .

# Expert Survey

"fmaximum\_numl"

1. fminimum\_magl
2. fminimuml
3. fminimum\_mag\_numl
4. fminimum\_numl

☐ Yes

☐ No

"execi"

1. j1l
2. exp10l
3. exp2l
4. expm1l

☐ Yes

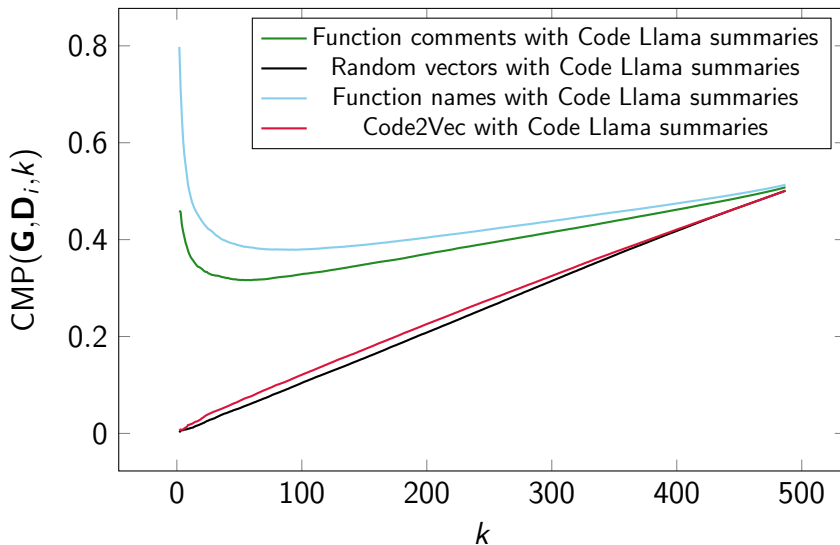
☐ No

Figure: Positive example

Figure: Negative example

Expert survey results				
Method	Code Llama summaries	Function names	Function comments	Code2Vec
Score	0.596	0.532	0.433	0.321

# Embedding space comparison



# Function names

Abbreviations can potentially confuse the SentenceTransformer:

Example function `lchmod`:

`l`  $\leftrightarrow$  link,    `ch`  $\leftrightarrow$  change,    `mod`  $\leftrightarrow$  file mode.

Nearest neighbors in function space:

`lchmod`  $\leftrightarrow$  (`lcong48`, `fchmodat`, `coshl`, `cacoshl`)

$\rightsquigarrow$  In categories:

`files`  $\leftrightarrow$  (`math`, `files`, `math`, `math`)

# function names

Example function `lchmod`:

`l` ↔ link,    `ch` ↔ change,    `mod` ↔ file mode.

Nearest neighbors in code llama summary space:

`lchmod`

↔ (`fchmodat`, `fchownat`, `euidaccess`, `__file_change_detection_for_stat`)

↪ In categories:

`files` ↔ (`files`, `files`, `files`, `files`)

# function comments

Comments are not always directly about the code:

Example functions `rand` and `rand_r`:

`rand` ↔ Return a random integer between 0 and `RAND_MAX`.

`rand_r` ↔ This algorithm is mentioned in the ISO C standard, here extended for 32 bits.

↪ Cosine distance in comment and llama summary space

$$d_{\text{comment}}(\text{rand}, \text{rand\_r}) = 0.8544 \quad d_{\text{llama}}(\text{rand}, \text{rand\_r}) = 0.2216.$$

# Future Work

## ► Code Llama

1. Is it necessary to use a large Model with 70B parameters?
2. Can Large Language Models produce deterministic output for this application?
3. Is there a better Prompt?

## ► Comments

1. Use inline Comments

# Conclusion

- ▶ Best strategies ranked:
  1. Code Llama summaries
  2. Function names
  3. Function comments
  4. Code2Vec
- ▶ Code Llama summary vectors for C source code downstream tasks
- ▶ Code Llama summary vectors can now be used to train a Model
- ▶  $\text{CMP}(A, B, k)$  function can be used to compare two embedding spaces from the same features Space
- ▶ Evaluation methods can be used to compare different Large Language Models to each other



# Embedding space comparison

$$\text{compare}(u, v)_k = \frac{1}{G_k} \sum_{i=1}^k \frac{\text{score}_k(u_i, i, v)}{\log_2(i+1)} \in [0, 1]$$

where

$u, v \in \mathbb{N}^k$  : Neighbor ranking of the same vector in different spaces,

$$\text{score}_k(l, i, v) = \begin{cases} 1 & , \exists j \in \mathbb{N} : l = v_j \wedge i = j \\ \frac{1}{2} & , \exists j \in \mathbb{N} : l = v_j \wedge i \neq j \\ 0 & , \text{otherwise} \end{cases} \quad , \quad G_k := \sum_{i=1}^k \frac{1}{\log_2(i+1)}.$$

# Embedding space comparison

$$\text{CMP}(A, B, k) = \frac{1}{N} \sum_{i=1}^N \text{compare}_k(\text{NN}_k(A_i, A), \text{NN}_k(B_i, B))$$

where

$A, B \in \mathbb{R}^{N \times l}$  : Embedding space with  $N$  vectors of length  $l$

$\text{NN}_k(A_i, A)$  :  $k$  nearest neighbors from vector with index  $i$  in  $A$

$k \in \mathbb{N}$  : Amount of vectors we include in one neighborhood relation

# Future Work

$$\text{CMP}(A, B, k) = \frac{1}{N} \sum_{i=1}^N \text{compare}_k(\text{NN}_k(A_i, A), \text{NN}_k(B_i, B))$$

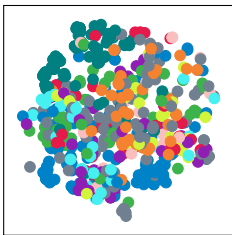
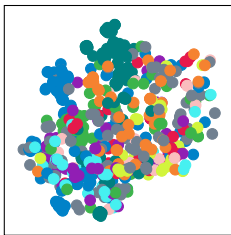
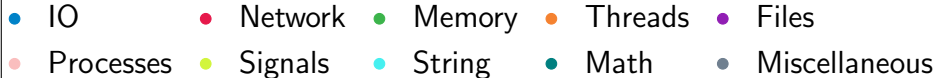
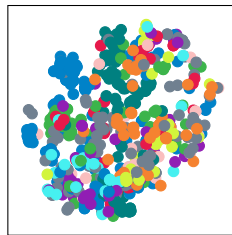
$$\text{compare}(u, v)_k = \frac{1}{G_k} \sum_{i=1}^k \frac{\text{score}_k(u_i, i, v)}{\log_2(i+1)} \in [0, 1]$$

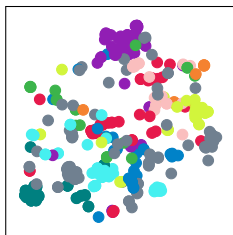
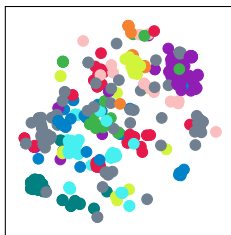
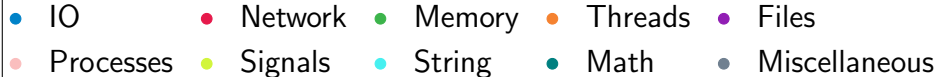
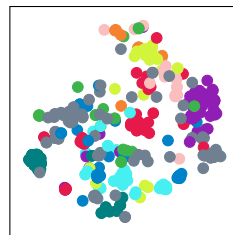
►  $\text{CMP}(A, B, k) \in [0, 1]$  function

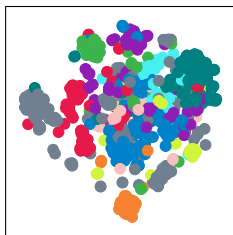
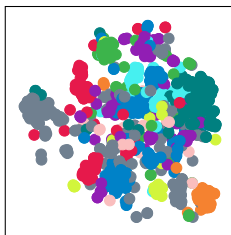
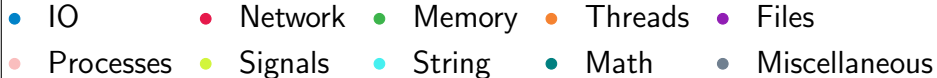
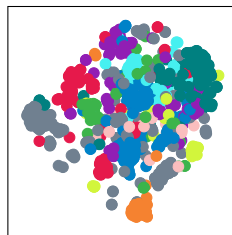
1. Is there an optimal value for  $k$ ?
2. Is there a better way to generate a neighborhood?  
(instead of K-Nearest-Neighbor)
3. Is there a better way to aggregate the compare functions?

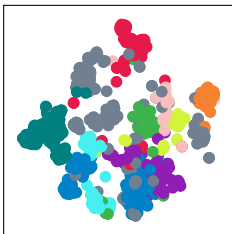
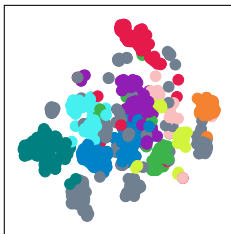
# Code2Vec

- ▶ Also dependent on the function names in the data set
- ▶ Bad results could be Explained by:
  1. Small data set
  2. C instead of Java  $\rightsquigarrow$  potential engineering mistakes
  3. Quality of names in the data set

Code2Vec  $P = 20$ Code2Vec  $P = 30$ Code2Vec  $P = 40$ 

Comments  $P = 20$ Comments  $P = 30$ Comments  $P = 40$ 

Names  $P = 20$ Names  $P = 30$ Names  $P = 40$ 

Code Llama  $P = 20$ Code Llama  $P = 30$ Code Llama  $P = 40$ 