



Comparing Natural Language Embeddings for Libc Functions as Rich Labels

Bachelor defense

Ruben Triwari

Ludwig Maximilian University Munich

19, February 2025

Outline

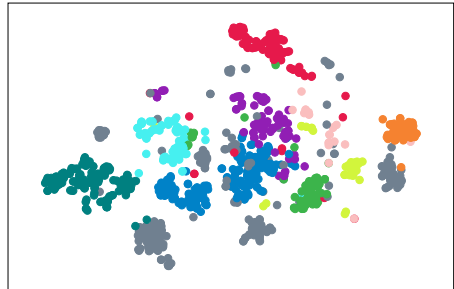
Motivation & Research Objective

Methodology

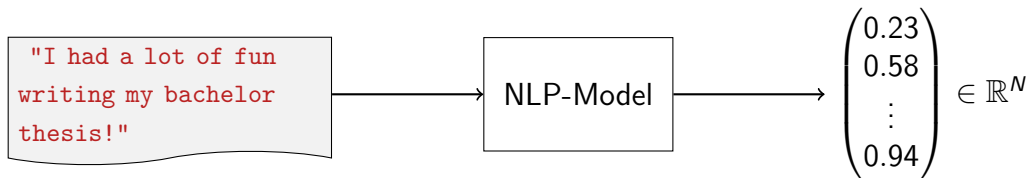
Results

Limitations

Conclusion & Future Work

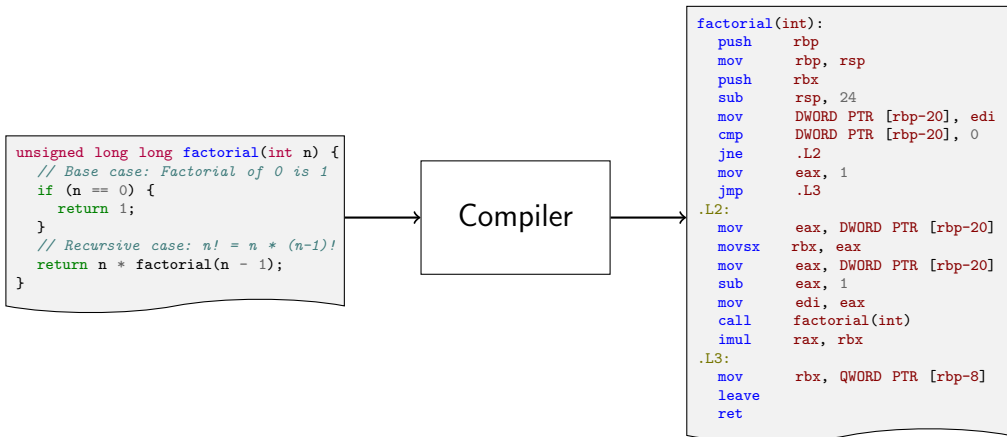


Motivation



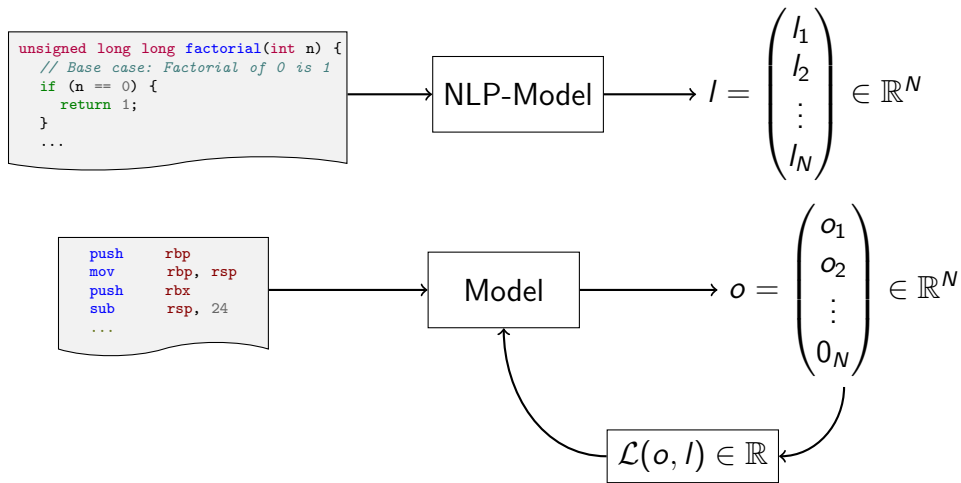
- ↪ Encoding natural language was a huge factor in recent nlp advancements
- ↪ Information described as a vector can be used in many downstream task
- ↪ That motivates encoding binary code and describing them as a vector
- ↪ That motivates using NLP tools to encode binary code

Motivation



⇒ Compiler removes important information in natural language

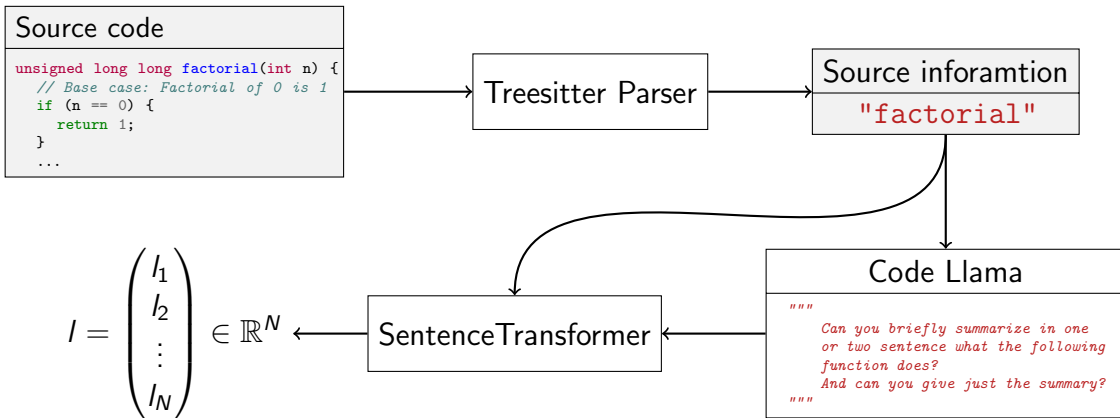
Motivation



Research Objectives

- ▶ Compare different 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 existing Code2Vec Model
- ▶ Propose a new way comparing embedding spaces

Architecture



Expert Survey

"execl"

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

☐ Yes

☐ No

"fmaximum_numl"

1. fminimum_magl
2. fminimuml
3. fminimum_mag_numl
4. fminimum_numl

☐ Yes

☐ No

Figure: Positive exmaple

Figure: Negative exmaple

Ergebnisse der Expertenbefragung				
Strategie	Code-Llama-Erklärungen	Funktionsnamen	Funktionskommentare	Code2Vec
Score	0.596	0.532	0.433	0.321

Embeddings 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,

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

Embeddings space comparison

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

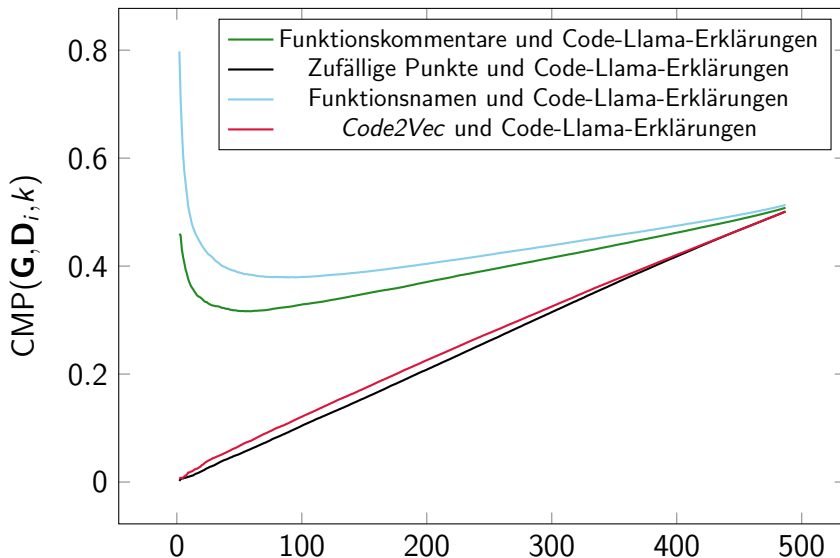
where

$A, B \in \mathbb{R}^{N \times l}$: 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

Embeddings space comparison



Evaluation with T-SNE

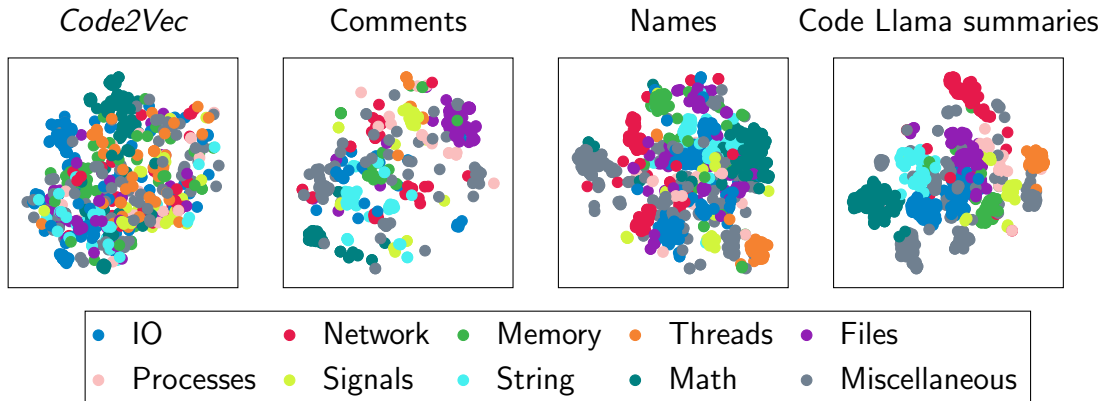


Figure: Depicted are the t -SNE output vectors with perplexity $P = 30$.

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

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

Code2Vec

- ▶ Also depend on the function names in the data set
- ▶ Bad results could be Explained by:
 1. Small data set
 2. C instead of Java
 3. Quality of names in the data set

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

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?

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

Discussion