

# Semantic Alignment Between Natural Language and Programming Language

Haochen Li HAOCHEN003@ntu.edu.sg AI-assisted Coding Changes How We Write Code



A coding assistant needs to understand requirements expressed in <u>natural</u> <u>language</u> and translate them into <u>programming language</u>.

#### Overview of My Research

#### One-to-One Mapping

– How can we align representations of natural language and programming language better, given a dataset that consists of NL-code pairs?

#### • One-to-Many Mapping

- In fact, an NL description often corresponds to multiple code implementations. Considering this, how can we further enhance semantic alignment?
- Many-to-Many Mapping
  - Similarly, a code snippet corresponds to multiple NL descriptions. Can we create highquality synthetic data by leveraging this relationship?

### **Representation Alignment**

#### Objective:

#### Similarity of a positive pair

Sort a given matrix in ascending order according to the sum of its rows.

def sort\_matrix(M):
 result = sorted(M, key=sum)
 return result

Sort a given matrix in ascending order according to the sum of its rows.

def square\_perimeter(a):
 perimeter=4\*a
 return perimeter

Similarity of any negative pairs

# Optimization with Contrastive Learning

Dominant choice: InfoNCE

$$\mathcal{L} = -\mathbb{E}_{i} \left[ \log \frac{\exp(text_{i} \cdot code_{i})}{\exp(text_{i} \cdot code_{i}) + \sum_{j \neq i} \exp(text_{i} \cdot code_{j})} \right]$$



Numerator: Pull together representations

Denominator: Push away representations

▲ Text ○ Positive Code ● Negative Code

### Augment Positive Pairs

Raw-data level augmentation methods are resource-consuming and limited in variety.



### Representation-level Augmentation

Existing Work:

Linear Interpolation e.g.  $0.8 \times h_1 + 0.2 \times h_2$ 

Stochastic Perturbation e.g.  $[h_{10}, h_{11}, h_{12}, h_{13}]$ -->  $[0, h_{11}, h_{12}, 0]$ 

New Method: Linear Extrapolation e.g.  $1.2 \times h - 0.2 \times h'$ 

> Binary Interpolation e.g.  $[h_{10}, h_{11}, h_{12}, h_{13}] \rightarrow [h_{10}, h'_{11}, h_{12}, h'_{13}]$

> > Gaussian Scaling e.g.  $h+\beta \odot h' \forall \beta \sim N(0, \sigma)$

General Format:

$$h^+ = \alpha \odot h + \beta \odot h'$$

Li, et al. "Exploring Representation-level Augmentation for Code Search." EMNLP 2022.

#### **Theoretical Analysis**

• Optimizing InfoNCE loss *L* improves the lower bound of mutual information *I* for a positive pair:

$$I \ge \log(B) - \mathcal{L}$$

• With augmentation, the lower bound is:  $I \ge \frac{1}{\alpha^2} (\log(NB) - \mathcal{L} - (\beta^2 + 2\alpha\beta) \cdot I^-)$ 

Moderate augmentation ( $\beta$  close to 0) leads to a **tighter** lower bound.

Mutual information indicates the similarity between two samples (the higher, the better)

Model	Ruby		JavaScript		Go		Python		Java		PHP	
	Original	w/ RA	Original	w/ RA	Original	w/ RA	Original	w/ RA	Original	w/ RA	Original	w/ RA
RoBERTa (code)	64.1	66.5	58.3	61.2	86.7	89.2	61.0	66.3	63.4	67.4	58.4	61.7
CodeBERT	64.8	66.4	59.4	60.8	87.8	<b>89.0</b>	63.6	65.4	66.3	67.4	61.5	61.9
GraphCodeBERT	70.5	72.1	64.7	67.1	89.6	90.3	69.0	70.8	69.1	70.8	64.8	65.6



Augmentation makes code representations more evenly distributed across the entire space.



Easier to distinguish representations.

Representation-level Augmentation also benefits Passage Retrieval



# Let's turn our sight to the negative pair side

We expect...

Similarity



#### Negative pairs should not be treated equally



How could we estimate  $\lambda_{ii}$ ? BM25, SimCSE, Fine-tuned Code Search Models...

Li, et al. "Rethinking Negative Pairs in Code Search." EMNLP 2023.

#### Soft-InfoNCE controls the distribution of negative samples

• **Theorem 1** For a batch of texts  $\{t_i\}_{i=1}^N$ , codes  $\{c_i\}_{i=1}^N$ , and similarity scores  $\{S_{ij}\}(S_{ij} \propto \lambda_{ij})$ , we have:

$$\mathcal{L}_{unif} \geq \frac{1}{N(N-2)} \sum_{i=1}^{N} \left[ \sum_{j\neq i}^{N} \log P_{\theta}(c_j|t_i) + KL(S_{ij}|P_{\theta}(c_j|t_i)) \right] + const.$$

Where *N* is the batch size,  $P_{\theta}(c_j | t_i)$  is predicted similarity between text  $t_i$  and code  $c_j$  by model  $\theta$ .

While  $L_{unif} \searrow$ , it encourages  $KL(S_{ij} | P_{\theta}(c_j | t_i)) \searrow$ the model predicted similarity  $P_{\theta}(c_j | t_i)$  gets closer to the given similarity.

#### Soft-InfoNCE reduces bias in mutual information estimation



Monte Carlo estimation for mutual information of all negative pairs  $C_{neg}$  based on the negative pairs in a batch

$$\mathbb{E}\log\left[1+\frac{p(c_i)}{p(c_i|t_i)}\sum_{j\neq i}^N \lambda_{ij}\frac{p(c_j|t_i)}{p(c_j)}\right]$$

By inserting  $\lambda_{ij}$ , it could be considered as using <u>importance sampling</u> to reduce estimation bias.

# Soft-InfoNCE outperforms InfoNCE

Model	Loss	Estimator	Ruby	Python	Java	JavaScript	PHP	Go
	InfoNCE	-	64.8	63.6	66.3	59.4	61.5	87.8
CodeBERT	Soft-InfoNCE	BM25 Trained Model SimCSE	66.0 <b>68.2</b> 66.6	66.4 <b>67.5</b> 66.8	<b>68.2</b> <b>68.2</b> 67.0	60.0 <b>61.5</b> 60.7	62.1 <b>63.1</b> 62.3	88.2 <b>89.8</b> 88.7
	InfoNCE	-	70.5	69.0	69.1	64.7	64.8	89.6
GraphCodeBERT	Soft-InfoNCE	BM25 Trained Model SimCSE	<b>73.0</b> 71.9 72.1	69.7 69.2 <b>70.0</b>	69.8 69.2 <b>70.2</b>	65.2 65.3 <b>65.6</b>	65.5 64.8 <b>65.7</b>	89.2 88.9 89.4
	InfoNCE	-	74.0	72.0	72.6	68.4	67.6	91.5
UniXCoder	Soft-InfoNCE	BM25 Trained Model SimCSE	75.3 75.3 75.3	<b>72.8</b> <b>72.8</b> 72.6	<b>73.3</b> 73.1 73.1	69.3 69.4 <b>69.9</b>	<b>68.4</b> 68.2 <b>68.4</b>	<b>91.6</b> 91.5 91.3

### Overview of My Research

#### One-to-One Mapping

– How can we align representations of natural language and programming language better, given a dataset that consists of NL-code pairs?

#### • One-to-Many Mapping

- In fact, an NL description often corresponds to multiple code implementations. Considering this, how can we further enhance semantic alignment?
- Many-to-Many Mapping
  - Similarly, a code snippet corresponds to multiple NL descriptions. Can we create highquality synthetic data by leveraging this relationship?

# Moving beyond One-to-One Mapping



In the first part, we discuss how to align representations better given NL-code pairs that are crawled from the Web (e.g. GitHub). There are many semantic-equivalent code variations. It is hard to collect all variations.

Ignoring this one-to-many relationship may not be optimal.

In the first part, we try to directly align natural language and programming language.

Some works translate natural language to programming language before calculating representation similarity.



Vanilla Code Search



Generation-Augmented Retrieval (GAR)

# GAR sometimes fails for code search

#### Text:

Write a function to get the frequency of the elements in a list.



\*Count the frequency manually

#### **Exemplar Code:**

```
def count frequency(my list):
    frequency = {}
    for element in my list:
        if element not in frequency:
            frequency[element] = 0
        frequency[element] += 1
    return frequency
```

\*Count the frequency automatically

True Code: import collections def freq count(list1): freq count =  $\setminus$ collections.Counter(list1) return freq\_count



LMs cannot recognize the similarity between the exemplar code and the true code due to their implementation style.

Li, et al. "Rewriting the Code: A Simple Method for Large Language Model Augmented Code Search." ACL 2024.

# Our Solution: rewrite the code (ReCo)



#### Framework Overview





Haochen Li

# Is the style mismatch problem in GAR a coincidence?

• Original code in the dataset is an implementation of the text in the dataset.  $code \sim P_{real}(\cdot | text)$ 

 $P_{real}$  refers to the real-world distribution.

• Consider the exemplar code generation process as a sampling process (LLM the sampler).

exemplar code~ $P_{LLM}(\cdot | \text{text})$ 

LLM also defines a probability distribution  $P_{LLM}$ .

Once there is a gap between  $P_{real}$  and  $P_{LLM}$ , the style mismatch is inevitable.

# Why ReCo works?

• Consider the exemplar code generation process as a sampling process (LLM the sampler).

exemplar code~ $P_{LLM}(\cdot | \text{text})$ LLM also defines a probability distribution  $P_{LLM}$ .

• We rewrite the code in the codebase through a code->summary->rewritten code process.

```
rewritten code~P_{LLM}(\cdot | \text{summary})
```

If the summary is perfect (summary  $\approx$  text): rewritten code  $\sim P_{LLM}(\cdot | \text{text})$ 

	CoNaLa	MBPP	APPS	MBJP		CoNaLa	MBPP	APPS	MBJP
Unsupervis	sed				Supervised				
BM25	52.6	12.6	11.6	11.3	CodeBERT	83.6	79.6	25.1	79.6
+ GAR	71.7	35.1	17.6	33.5	+ GAR	88.6	87.7	<u>    29.3                                </u>	84.1
+ ReCo	<b>75.8</b> <sub>+4.1</sub>	<b>70.8</b> +35.7	$22.6_{+5.0}$	<b>65.3</b> <sub>+31.8</sub>	+ ReCo	$85.0_{-3.6}$	<b>92.3</b> <sub>+4.6</sub>	<b>51.2</b> <sub>+21.9</sub>	$89.1_{+5.0}$
UniXcoder	77.2	69.3	8.3	73.2	UniXcoder	84.8	81.2	24.3	81.6
+ GAR	83.9	85.0	13.2	80.0	+ GAR	85.9	89.0	<u>34.5</u>	85.6
+ ReCo	$85.1_{+1.2}$	<b>92.4</b> $_{+7.4}$	<b>28.8</b> <sub>+15.6</sub>	<b>87.6</b> <sub>+7.6</sub>	+ ReCo	$87.1_{+1.2}$	$94.2_{+5.2}$	$58.1_{+23.6}$	$90.5_{+4.9}$
Contriever	55.7	55.3	9.6	37.0					
+ GAR	75.0	71.3	_14.0	<u>62.3</u>					
+ ReCo	$77.9_{+2.9}$	<b>87.4</b> <sub>+16.1</sub>	<b>41.6</b> +27.6	$76.6_{+14.3}$	CAR in	deed in	nroves	code sea	rch
CodeT5+	73.7	59.4	7.6	67.7	but ReC	'o impr	$n_{\rm MP}$	re	
+ GAR	80.3	7 <u></u>	10.2	79.2				JIC.	<b>1a</b>
+ ReCo	<b>80.8</b> <sub>+0.5</sub>	<b>89.4</b> <sub>+11.7</sub>	<b>29.9</b> <sub>+19.7</sub>	<b>84.0</b> <sub>+4.8</sub>	(non-ne tune)	ural mo	odel, zei	ro-shot, f	ine-
+ KeCo	<b>ð0.ð</b> +0.5	<b>89.4</b> +11.7	<b>29.9</b> +19.7	<b>84.0</b> +4.8	tune)		Juci, Zci	10-31101, 1	111

	CoNaLa	MBPP	APPS	MBJP
Unsupervis	red			
BM25	52.6	12.6	11.6	11.3
+ GAR	71.7	35.1	17.6	33.5
+ ReCo	<b>75.8</b> <sub>+4.1</sub>	<b>70.8</b> +35.7	<b>22.6</b> +5.0	<b>65.3</b> +31.8
UniXcoder	77.2	69.3	8.3	73.2
+ GAR	83.9	85.0	13.2	80.0
+ ReCo	$85.1_{+1.2}$	<b>92.4</b> <sub>+7.4</sub>	$28.8_{+15.6}$	<b>87.6</b> <sub>+7.6</sub>
Contriever	55.7	55.3	9.6	37.0
+ GAR	75.0	71.3	14.0	62.3
+ ReCo	$77.9_{+2.9}$	<b>87.4</b> <sub>+16.1</sub>	$41.6_{+27.6}$	<b>76.6</b> <sub>+14.3</sub>
CodeT5+	73.7	59.4	7.6	67.7
+ GAR	80.3	77.7	10.2	79.2
+ ReCo	$80.8_{+0.5}$	<b>89.4</b> <sub>+11.7</sub>	$29.9_{+19.7}$	$84.0_{+4.8}$

	CoNaLa	MBPP	APPS	MBJP
Supervised				
CodeBERT	83.6	79.6	25.1	79.6
+ GAR	88.6	87.7	29.3	84.1
+ ReCo	$85.0_{-3.6}$	<b>92.3</b> <sub>+4.6</sub>	<b>51.2</b> <sub>+21.9</sub>	$89.1_{+5.0}$
UniXcoder	84.8	81.2	24.3	81.6
+ GAR	85.9	89.0	34.5	85.6
+ ReCo	$87.1_{+1.2}$	$94.2_{+5.2}$	$58.1_{+23.6}$	$90.5_{+4.9}$

With ReCo, non-neural model BM25 is comparable to neural models under the zero-shot setting.

	CoNaLa	MBPP	APPS	MBJP		CoNaLa	MBPP	APPS	MBJP
Unsupervis	red								
BM25	52.6	12.6	11.6	11.3	CodeBERT	83.6	79.6	25.1	79.6
+ GAR	71.7	35.1	17.6	33.5	+ GAR	88.6	87.7	29.3	84.1
+ ReCo	$75.8_{+4.1}$	$70.8_{+35.7}$	$22.6_{+5.0}$	<b>65.3</b> <sub>+31.8</sub>	+ ReCo	$85.0_{-3.6}$	<b>92.3</b> <sub>+4.6</sub>	<b>51.2</b> <sub>+21.9</sub>	$89.1_{+5.0}$
UniXcoder	77.2	69.3	8.3	73.2	UniXcoder	84.8	81.2	24.3	81.6
+ GAR	83.9	85.0	13.2	80.0	+ GAR	85.9	89.0	34.5	85.6
+ ReCo	<b>85.1</b> <sub>+1.2</sub>	<b>92.4</b> <sub>+7.4</sub>	$28.8_{+15.6}$	<b>87.6</b> <sub>+7.6</sub>	+ ReCo	$87.1_{+1.2}$	$94.2_{+5.2}$	$58.1_{+23.6}$	$90.5_{+4.9}$
Contriever	55.7	55.3	9.6	37.0					
+ GAR	75.0	71.3	14.0	62.3					
+ ReCo	$77.9_{+2.9}$	$87.4_{+16.1}$	$\textbf{41.6}_{+27.6}$	<b>76.6</b> +14.3	With De		unal max	dolound	or the
CodeT5+	73.7	59.4	7.6	67.7	zero-shot setting is comparable to them after fine-tuning.				
+ GAR	80.3	77.7	10.2	79.2					
+ ReCo	$80.8_{+0.5}$	<b>89.4</b> <sub>+11.7</sub>	<b>29.9</b> <sub>+19.7</sub>	$84.0_{+4.8}$					

#### A new metric to measure code style similarity

- Code Style Similarity considers similarity of style from three perspectives:
  - Variable Naming (based on Edit Distance)
  - API Invocation (based on Edit Distance)
  - Code Structure (based on Tree Edit Distance of Abstract Syntax Trees)

$$\operatorname{CSSim}(c_1, c_2) = 1 - \frac{1}{3} (\operatorname{Dis}_{\operatorname{Var}} + \operatorname{Dis}_{\operatorname{API}} + \operatorname{TED})$$
  
$$\operatorname{Dis} = \frac{1}{2} (\frac{1}{||\lambda||_1} \sum_{v_i \in \mathbf{S}_1} \lambda_i \min_{v_j \in \mathbf{S}_2} \operatorname{ED}(v_i, v_j) + \frac{1}{||\lambda||_1} \sum_{v_i \in \mathbf{S}_2} \lambda_i \min_{v_j \in \mathbf{S}_1} \operatorname{ED}(v_i, v_j))$$

*ED*: Edit Distance *TED*: Tree Edit Distance of Abstract Syntax Tree  $\lambda_i$ : Normalized Inverse Document Frequency (IDF)  $S_1, S_2$ : Set of extracted variables or APIs

### Our metric is better than existing metrics



The improvement of Code Style Similarity has stronger correlation with the performance improvement of ReCo over GAR.

# Our metric is better than existing metrics

Dataset	Random	w/ the best exemplar code						
Dutuset		CSSim	CodeBLEU	ROUGE-L	BLEU			
CoNaLa	41.3	43.6	39.6	41.5	42.3			
MBPP	24.0	26.8	25.1	23.7	24.0			
APPS	14.1	15.2	14.8	14.2	14.2			
MBJP	26.6	29.9	29.1	27.2	28.8			

We select the best exemplar code (the most similar to true code) under different metrics for BM25+ GAR.

The Code Style Similarity outperforms other settings.

True Code	Exemplar Code 1	Exemplar Code 2
<pre>def find_first_duplicate(nums):     num_set = set()     no_duplicate = -1     for i in range(len(nums)):         if nums[i] in num_set:             return nums[i]         else:             num_set.add(nums[i])     return no_duplicate</pre>	<pre>def find_duplicate(nums):     for i in range(len(nums)):         if nums[i] in nums[i+1:]:             return nums[i]     return None         Preferred by         preferred by         other Metrics         other</pre>	<pre>def find_duplicate(my_list):     seen = set()     for num in my_list:         if num in seen:             return num         seen.add(num)     return None         Preferred by         Preferred by         OMr Metric         OMr Metric </pre>

### Overview of My Research

#### One-to-One Mapping

- How can we align representations of natural language and programming language better, given a dataset that consists of NL-code pairs?

#### • One-to-Many Mapping

– In fact, an NL description often corresponds to multiple code implementations. Considering this, how can we further enhance semantic alignment?

#### • Many-to-Many Mapping

– Similarly, a code snippet corresponds to multiple NL descriptions. Can we create highquality synthetic data by leveraging this relationship?

# **Iterative Rejection Fine-Tuning**



# There is bias in data synthesis process

For each seed input, it is a specific expression of the <u>intent behind the input.</u>

Intent:

– I want a function for substring matching.

Possible inputs 1:

– Write a function to search a string for a regex pattern.

Possible input 2:

– Find a substring within a larger string that matches a given regular expression pattern.

# There is bias in data synthesis process

For each seed input, it is a specific expression of the <u>intent behind the input</u>. For each intent, suppose all possible inputs and valid codes form a joint space.



If we synthesize codes solely based on each seed input  $d_0$ , it is drawing codes from a **conditional distribution** (Red).

It is better to sample codes from the **marginal distribution** (Blue).

Li, et al. "GiFT: Gibbs Fine-Tuning for Code Generation." ACL 2025.

## Why is marginal distribution better?

• The loss  $\mathcal{L}_{marg}$  is implicitly estimated over all possible  $\mathcal{L}_{cond}$ .

$$\mathcal{L}_{marg} = -\mathbb{E}_{d \sim P(input)} \mathcal{L}_{cond}(d)$$

• The variance (diversity) of code is higher.

$$Var(c) = \mathbb{E}_d[Var(c|d)] + Var(\mathbb{E}_d[c|d])$$

Expected variance (diversity) of codes from conditional distribution

Variance of codes from different inputs

# How can we sample from marginal distribution?

- One possible solution is to ask LLMs to rewrite the seed input.
  - -The quality of rewritten inputs is not satisfiable.
  - -Rewritten inputs are still similar.

Datasets		GiFT	
	Seed	Rewritten	
APPS+ (Intro.)	17.79	4.8	>10.31
APPS+ (Inter.)	3.22	0.38	>2.28
MBPP+	53.71	24.6	>24.07
CodeInsight	43.87	9.84	>24.18

Pass rate (%) of self-generated codes from the seed description, rewritten description, and GiFT.

• We got inspiration from Gibbs Sampling, commonly used to approximate joint distributions based on conditional distributions.

# Introduction of Gibbs Sampling

Take a two-variable joint distribution as an example. We want to sample from P(x, y), but we only have access to P(x|y) and P(y|x).

```
Starting from x_0
For t = 1 to N do:
y_t \sim P(y|x_{t-1})
x_t \sim P(x|y_t)
```

Collected pairs { $(x_1, y_1)(x_2, y_2) \dots (x_N, y_N)$ } can be considered as being sampled from the joint distribution.

# Introduction of Gibbs Sampling

#### **Gibbs sampling**

We want to sample from P(x, y), but we only have access to P(x|y) and P(y|x).

Starting from  $x_0$ For t = 1 to N do:  $y_t \sim P(y|x_{t-1})$  $x_t \sim P(x|y_t)$ 

#### In our context

We want to sample from *P*(text, code). We have access to: *P*(text|code) ---- Code Summarization *P*(code|text) ---- Code Generation

Starting from the seed input text<sub>0</sub>

Collected pairs  $\{(text_1, code_1) \dots (text_N, code_N)\}$ 

# Curating codes for fine-tuning



The marginal distribution (Blue) often does not follow a uniform distribution. Using all correct codes for fine-tuning:

- 1. Bias towards easier inputs
- 2. Bias towards head in long-tail distribution

#### Solutions:

- 1. At most K codes per inputs
- 2. Weighted random sampling based on perplexity to encourage selecting tail codes



RFT refers to Rejection-Fine-tuning, RFT+RD refers to RFT with rewritten texts

Haochen Li

# Impact of Perplexity-guided Data Selection



T > 0: Encourage Tail T < 0: Encourage Head Encouraging head speeds up convergence, yet its potential is lower than encouraging tail.

Perplexity distribution of selfgenerated codes.

Encouraging head makes the distribution sharper.

### GiFT benefits not only self-training, but also distillation



Distill DeepSeek-Coder-6.7B to CodeLlama-7B

# Prerequisites of applying GiFT to other tasks?

- LLMs could seamlessly translate between inputs and outputs.
  - Counter example: Math Reasoning
  - Translate solutions back to math problems?
- There is indeed a joint input-output space.
  - Counter example: Factual Question Answering



Thank you! Any Questions?