

# Exploring Representation-Level Augmentation for Code Search

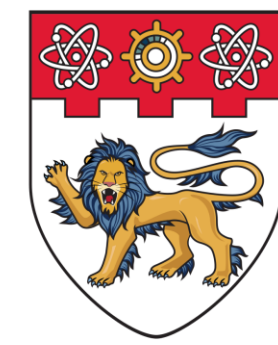
Haochen Li<sup>1</sup>, Chunyan Miao<sup>1,2</sup>, Cyril Leung<sup>1,2</sup>, Yanxian Huang<sup>3</sup>, Yuan Huang<sup>3</sup>, Hongyu Zhang<sup>4</sup>, and Yanlin Wang<sup>3</sup>

<sup>1</sup>Nanyang Technological University, Singapore

<sup>2</sup>China-Singapore International Joint Research Institute

<sup>3</sup>Sun Yat-sen University <sup>4</sup>The University of Newcastle

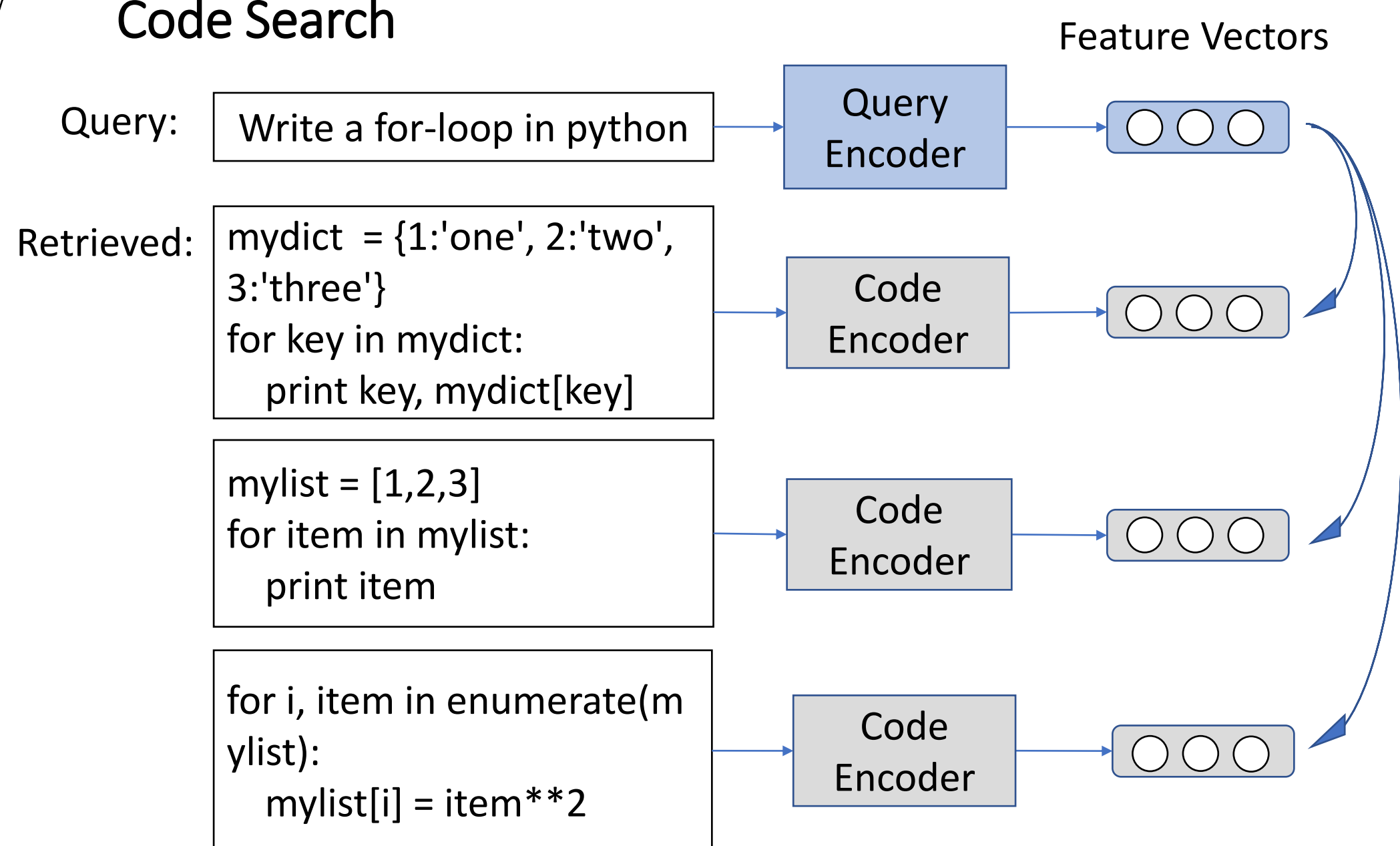
Paper: <https://arxiv.org/abs/2210.12285>



**NANYANG  
TECHNOLOGICAL  
UNIVERSITY**  
SINGAPORE



## Code Search



Code search, which aims at retrieving the most relevant code fragment for a given natural language query, is a common activity in software development practice.

## Theoretical Analysis

**Theorem 1** Optimizing InfoNCE loss  $L_N$  improves lower bounds of mutual information  $I(q, c)$  for a positive pair:

$$I(q, c) \geq \log(B) - L_N$$

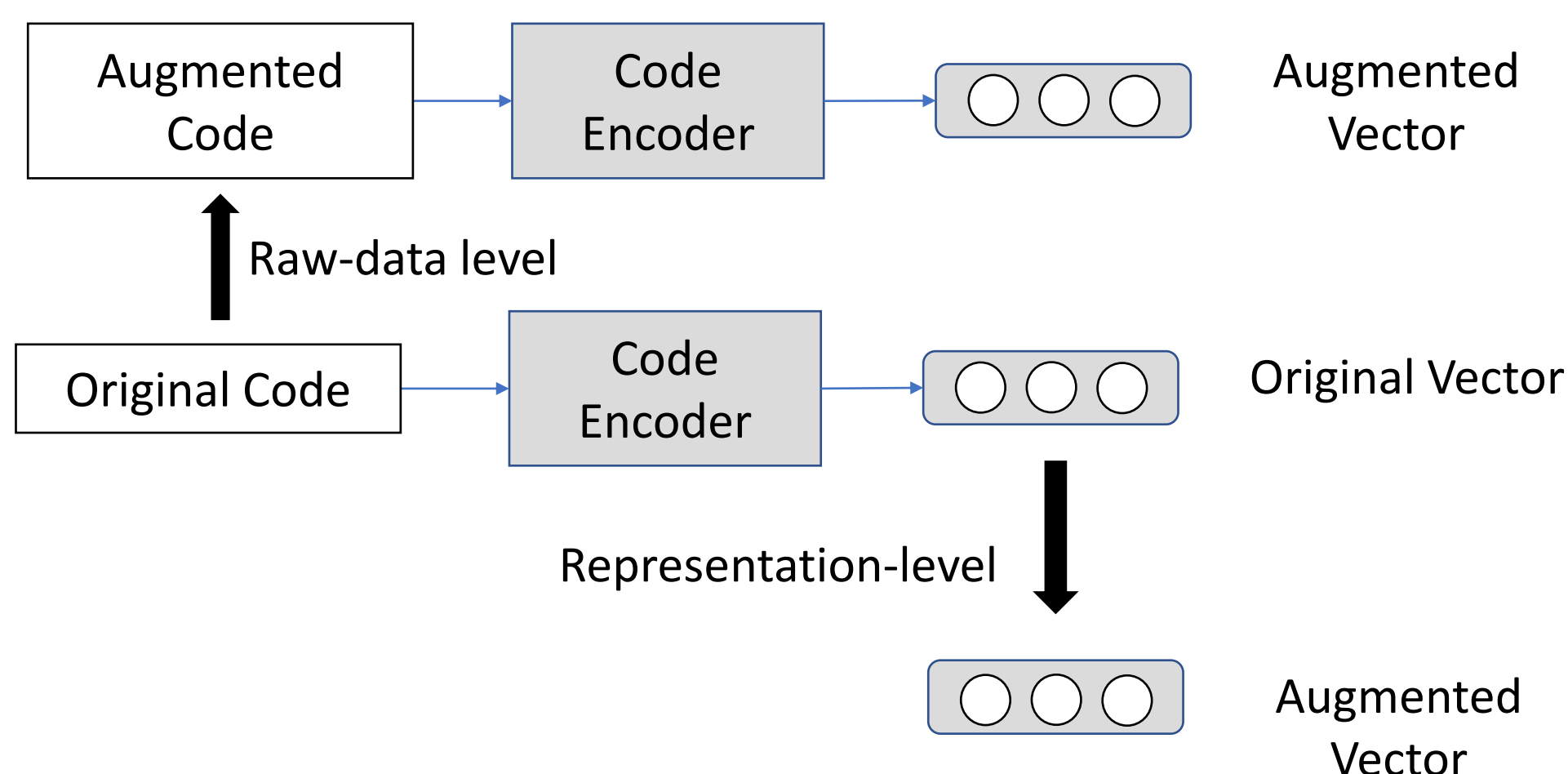
where  $q \in Q, c \in C$ , and  $B$  is the size of sets.

**Theorem 2** Optimizing InfoNCE loss  $L_N$  with representation-level augmentation improves a **tighter** lower bounds of mutual information  $I(q, c)$  for a positive pair:

$$I(q, c) \geq \frac{1}{\alpha^2} (\log(NB) - L_N - \alpha\beta \cdot I(q, c^-) - \alpha\beta \cdot I(q^-, c) - \beta^2 \cdot I(q^-, c^-))$$

where  $q, q^- \in Q, c, c^- \in C, (q, c^-), (q^-, c)$  and  $(q^-, c^-)$  are all negative pairs,  $\alpha$  and  $\beta$  are coefficients,  $B$  is the size of sets, and  $N$  is the augmentation time.

## Motivation

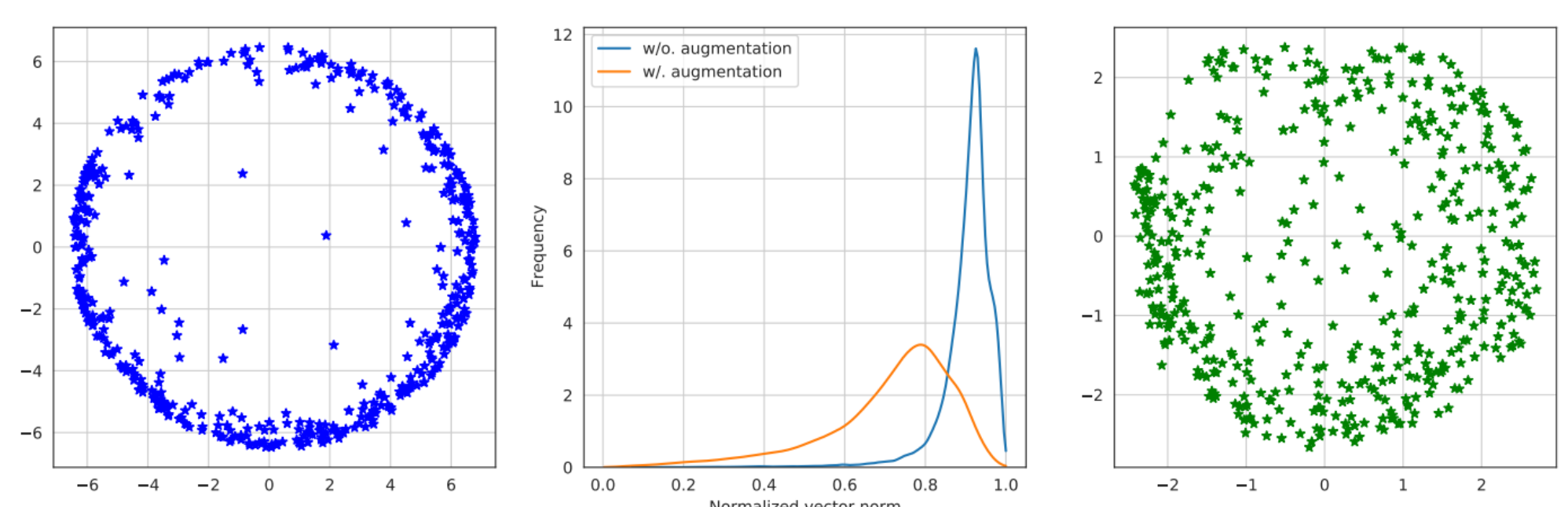


Existing augmentation methods are applied to raw data, which is **resource-consuming** and **limited**. For representation-level augmentation, we do not need to encode raw data again.

## Experiments

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)	0.641	<b>0.665</b>	0.583	<b>0.612</b>	0.867	<b>0.892</b>	0.610	<b>0.663</b>	0.634	<b>0.674</b>	0.584	<b>0.617</b>
CodeBERT	0.648	<b>0.664</b>	0.594	<b>0.608</b>	0.878	<b>0.890</b>	0.636	<b>0.654</b>	0.663	<b>0.674</b>	0.615	<b>0.619</b>
GraphCodeBERT	0.705	<b>0.721</b>	0.647	<b>0.671</b>	0.896	<b>0.903</b>	0.690	<b>0.708</b>	0.691	<b>0.708</b>	0.648	<b>0.656</b>

Results on CodeSearchNet Dataset. Performance of different models under MRR. "w/ RA" stands for "with representation-level augmentation".



Visualization of code vector distribution with and without representation-level augmentation. **Left:** without augmentation. **Right:** with augmentation. **Middle:** distribution of vector norms. With augmentation, representations better take advantage of norms.

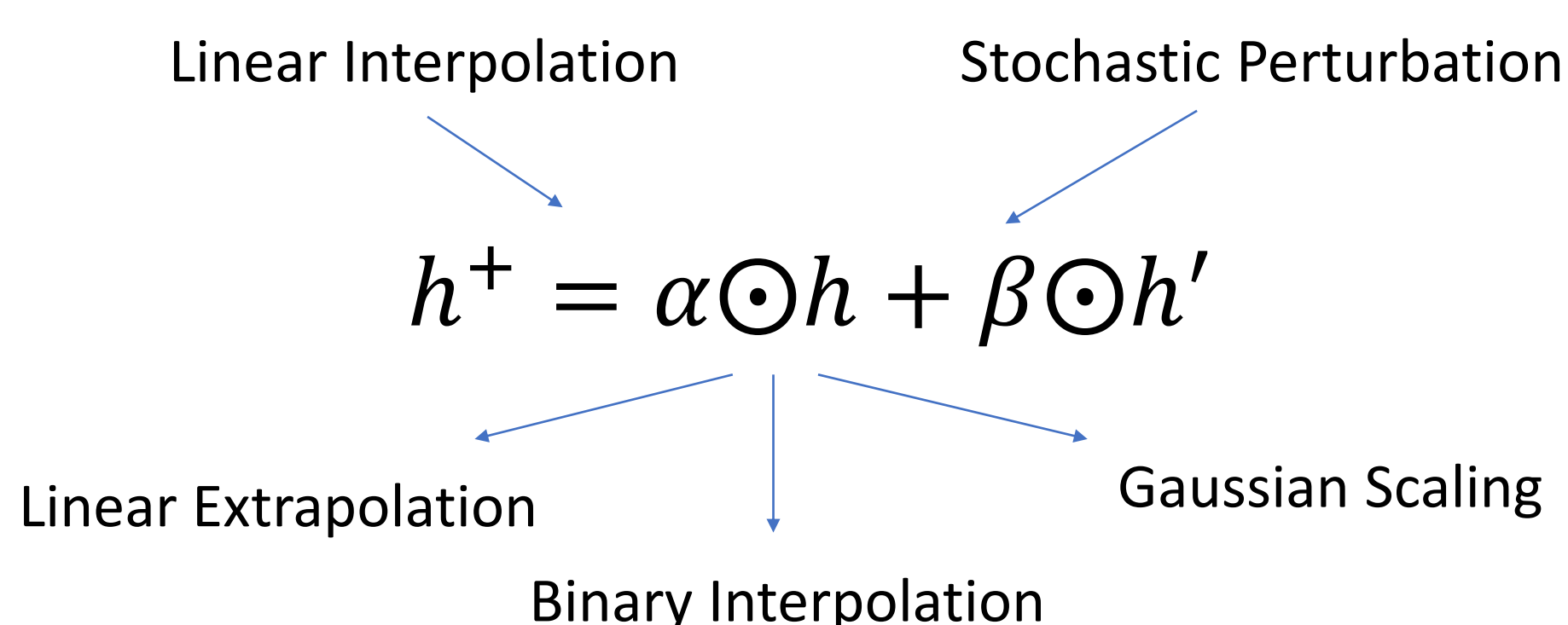
Augmentations	RoBERTa	CodeBERT	GraphCodeBERT
no augmentations	0.629	0.636	0.690
linear interpolation	0.644	0.648	0.702
linear extrapolation	0.640	0.646	0.704
stochastic perturbation	0.658	0.648	0.698
binary interpolation	0.655	<b>0.655</b>	0.705
Gaussian scaling	0.657	0.649	0.696
all augmentations	<b>0.663</b>	0.654	<b>0.708</b>

Five augmentation methods have **similar** effects on MRR since they come from the same general format.

Model	FiQA-2018		NFCorpus	
	Original	w/ RA	Original	w/ RA
DistilBERT	0.352	<b>0.400</b>	0.481	<b>0.505</b>
RoBERTa	0.343	<b>0.356</b>	0.367	<b>0.389</b>

We also evaluate the proposed methods on two **passage retrieval** datasets.

## Methodology



We unify the existing two approaches, **Linear Interpolation** and **Stochastic Perturbation**, to a general format. Based on it, we propose three novel augmentation approaches, **Linear Extrapolation**, **Binary Interpolation**, and **Gaussian Scaling**.

• Code: <https://github.com/Alex-HaochenLi/RACS>

• Contact: [haochen003@ntu.edu.sg](mailto:haochen003@ntu.edu.sg)

