

# CoCoSoDa: Effective Contrastive Learning for Code Search

Ensheng Shi<sup>a,†</sup> Yanlin Wang<sup>b,§</sup> Wenchao Gu<sup>c,†</sup> Lun Du<sup>d</sup>

Hongyu Zhang<sup>e</sup> Shi Han<sup>d</sup> Dongmei Zhang<sup>d</sup> Hongbin Sun<sup>a,§</sup>

<sup>a</sup>Xi'an Jiaotong University <sup>b</sup>School of Software Engineering, Sun Yat-sen University

<sup>c</sup>The Chinese University of Hong Kong <sup>d</sup>Microsoft Research <sup>e</sup>Chongqing University

s1530129650@stu.xjtu.edu.cn, wangyin36@mail.sysu.edu.cn

wcgu@cse.cuhk.edu.hk, {lun.du, shihan, dongmeiz}@microsoft.com

hyzhang@cqu.edu.cn, hsun@mail.xjtu.edu.cn

**Abstract**—Code search aims to retrieve semantically relevant code snippets for a given natural language query. Recently, many approaches employing contrastive learning have shown promising results on code representation learning and greatly improved the performance of code search. However, there is still a lot of room for improvement in using contrastive learning for code search. In this paper, we propose CoCoSoDa to effectively utilize contrastive learning for code search via two key factors in contrastive learning: data augmentation and negative samples. Specifically, soft data augmentation is to dynamically masking or replacing some tokens with their types for input sequences to generate positive samples. Momentum mechanism is used to generate large and consistent representations of negative samples in a mini-batch through maintaining a queue and a momentum encoder. In addition, multimodal contrastive learning is used to pull together representations of code-query pairs and push apart the unpaired code snippets and queries. We conduct extensive experiments to evaluate the effectiveness of our approach on a large-scale dataset with six programming languages. Experimental results show that: (1) CoCoSoDa outperforms 18 baselines and especially exceeds CodeBERT, GraphCodeBERT, and UniXcoder by 13.3%, 10.5%, and 5.9% on average MRR scores, respectively. (2) The ablation studies show the effectiveness of each component of our approach. (3) We adapt our techniques to several different pre-trained models such as RoBERTa, CodeBERT, and GraphCodeBERT and observe a significant boost in their performance in code search. (4) Our model performs robustly under different hyper-parameters. Furthermore, we perform qualitative and quantitative analyses to explore reasons behind the good performance of our model.

**Index Terms**—code search, contrastive learning, soft data augmentation, momentum mechanism

## I. INTRODUCTION

Code search plays an important role in software development and maintenance [1], [2]. To implement a certain functionality, developers often search and reuse previously-written code from open source repositories such as GitHub or from a large local codebase [3]–[5]. The key challenge in this task is to find semantically relevant code snippets *written in programming languages* based on input queries *written in natural languages* [6].

Early studies [3], [7]–[9] on code search mainly leverage on the lexical information of the code snippets and use information retrieval (IR) methods. Later on, deep learning-based approaches [10]–[21] that employ neural networks to embed code and queries into a shared embedding space and measure their semantic similarity through vector distances are explored. Recently, large pre-trained code models [22]–[26], which are pre-trained on large multi-programming-language data, improve the understanding of code semantics and achieve better code search performance. Some studies [27], [28] apply contrastive learning to unsupervised code representation learning and Corder [28] also achieves good performance in code search. In detail, they first use semantic-preserving transformations [28] to generate similar code snippets as positive samples. These transformations are including *Statements Permutation* [28] (swapping two statements that have no data dependency on each other in a basic block), *Loop Exchange* [28] (replacing *for* statements with *while* statements or vice versa.), etc. Second, they treat other code snippets in the same mini-batch as negative samples and then optimize the model to pull together representations of positive samples and push apart representations of negative samples. These approaches have shown promising results in code search. For example, Corder [28] outperforms all approaches and achieve a MRR score of 0.727 on CodeSearchNet [29] Java dataset. However, there is still a lot of room for improvement in leveraging contrastive learning for code search, such as using other effective data augmentation to generate positive samples or enriching negative samples.

In this paper, we propose **CoCoSoDa** (stands for Code search with multimodal momentum Contrastive learning and Soft Data augmentation) to effectively utilize contrastive learning for code search via two key factors in contrastive learning: data augmentation and negative samples. ① To learn a better overall semantic representation of the code snippet and query instead of focusing on token-level semantic representation learning according to the surrounding context, we propose four SoDa (stands for Soft Data augmentation) methods (Sec. III-C). They dynamically replace  $r$  ( $r$  is stable from 5% to 20%) of code tokens with their types or simply mask

<sup>§</sup>Yanlin Wang and Hongbin Sun are the corresponding authors.

<sup>†</sup>Work done during the author’s employment at Microsoft Research Asia

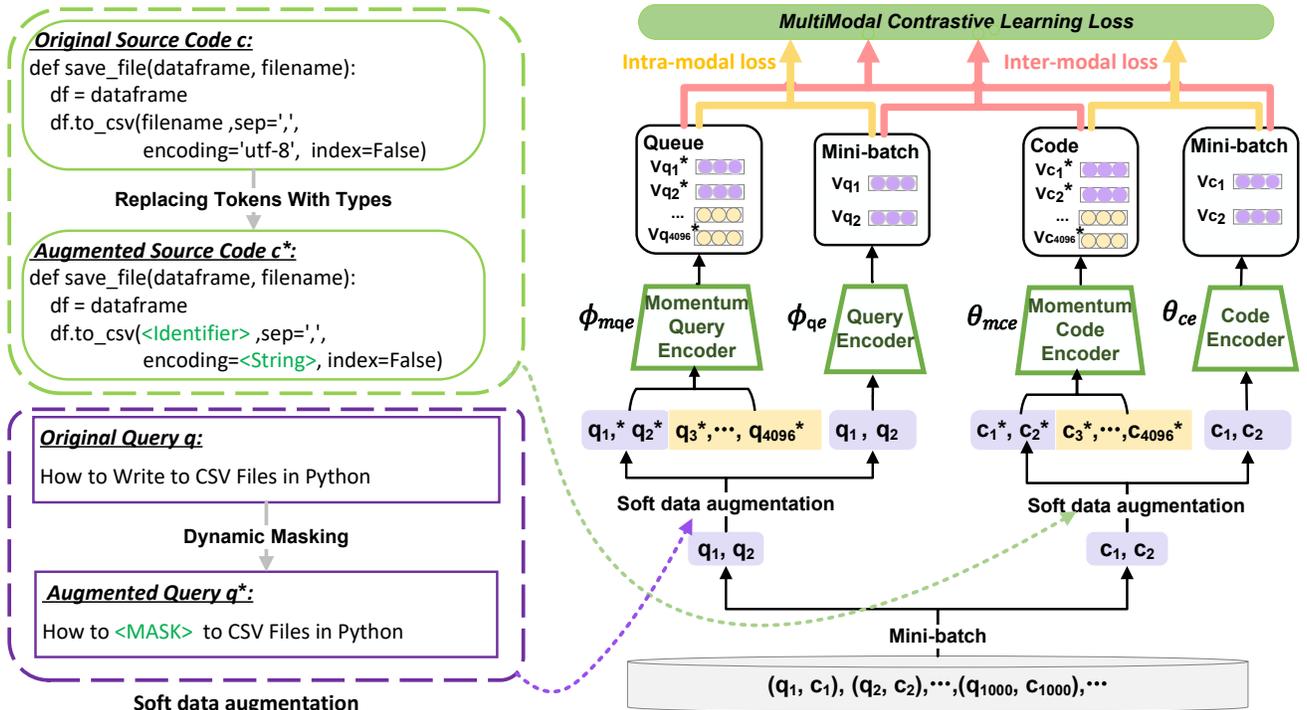


Fig. 1. The framework of CoCoSoDa.

them to generate similar code snippets. ② To distinguish one sample from more negative samples at each iteration, we adopt the momentum mechanism [30] to enlarge negative samples in a mini-batch. ③ We also employ multimodal contrastive learning to minimize the distance between the representations of code-query pair and maximize the distance between the representations of the query (code snippet) and other many unpaired code snippets (queries). The overall framework of CoCoSoDa is shown in Fig. 1. On the left is an example of soft data augmentation, and on the right is the main architecture of our model. More details are given in Sec. III.

We evaluate the effectiveness of CoCoSoDa on a large-scale dataset CodeSearchNet [29] with six programming language (Ruby, JavaScript, Go, Python, Java, PHP) and compare CoCoSoDa with 18 state-of-the-art (SOTA) approaches. We also conduct the ablation study to study the effectiveness of each component of CoCoSoDa. Furthermore, we apply CoCoSoDa to other three large-scale pre-trained models, including natural language pre-trained model RoBERTa [31], code pre-trained models CodeBERT [23] and GraphCodeBERT [22]. We also assign different hyperparameters to check their impact on code search. In addition, we discuss why CoCoSoDa perform well through qualitative and quantitative analyses. Experimental results show that: (1) CoCoSoDa significantly outperforms existing SOTA approaches on code search task (Sec. V-A). (2) The multimodal momentum contrastive learning including intra-modal and inter-modal contrastive learning and four SoDa methods play important roles individually and can improve the performance of the code search model (Sec. V-B). (3) CoCoSoDa can be easily adapted to other pre-trained

models and obviously boost their performance (Sec. V-C). (4) CoCoSoDa performs stably over a range of hyperparameters: learning rate is from  $5e^{-6}$  to  $7e^{-5}$ , momentum coefficient  $m$  is between 0.910 and 0.999, masked ratio  $r$  is from 5% to 20%, and temperature hyperparameter  $\tau$  varies from 0.03 to 0.07 (Sec. V-D).

We summarize the contributions of this paper as follows:

- We adapt Transform-based momentum contrastive learning algorithm to better leveraging contrastive learning techniques for code search task. It enables the model to learn effective code representations by learning better representation of one sample against more negative samples. We also propose a new approach incorporating multimodal momentum contrastive learning. It can pull together the representations of matched code-query pairs and push apart the representations of unmatched code-query pairs.
- We propose four simple yet effective soft data augmentation methods that utilize dynamic masking and replacement for data augmentation. More importantly, SoDa can be easily applied to all programming languages.
- We conduct extensive experiments to evaluate the superiority of our approach on a large-scale multi-programming-language dataset. The results show that our approach significantly outperforms baselines, our framework can be easily adapted to other pre-trained models and significantly boost their performance, and CoCoSoDa performs stably over a range of hyperparameters.

---

## II. RELATED WORK

### A. Code Search

Learning the representation of code is an emerging topic and has been found to be useful in many software engineering tasks, such as code summarization [32]–[36], code search [10], [15], [18], [37], [38], code completion [39]–[43], commit message generation [44]–[48]. Among them, code search plays an important role in software development and maintenance [1], [2]. Traditional approaches [3], [7]–[9] based on retrieval information techniques mainly focus on the lexical information of the source code and apply keywords matching methods to search the relevant code snippets for the given query. In recent years, deep learning-based approaches leverage the neural network to learn the semantic representations of the source code and natural language to improve the understanding of code snippets and queries. Gu et al. [10] is the first to use the deep neural network to embed the code and query into a shared vector space and measure the similarity of them using vector distance. Subsequently, various types of model structures are applied to code search, including sequential models [16]–[20], convolutional neural network [14], [15], [21], tree neural network [20], graph models [13], [20], and transformers [11], [15]. Recently, large-scale code pre-trained models [22]–[26], [49], which are pre-trained on a massive source code dataset, improve the understanding of code semantics and achieve significant improvements in code search task. For example, CodeBERT is pre-trained with masked language modelling (MLM), which is to predict masked tokens, and replaced token detection (RTD), which uses a discriminator to identify replaced tokens. GraphCodeBERT takes source code, paired summarization and the corresponding data flow as the input and is pre-trained with MLM, data flow edge prediction, and node alignment tasks. PLBART [26] is a sequence-to-sequence code pre-trained models and is pre-trained with denoising autoencoding, which destroys a span of tokens and then recovers them. Our approach can be easily applied to these pre-trained models and boost their performance.

### B. Code Representation Learning with Contrastive Learning

Contrastive learning approaches [50], which pull close the similar representations and push apart different representations, have been successfully used in self-supervised representation learning on images [30], [51] and natural language texts [52]–[54]. To generate individual augmentations, images usually use spatial [55], [56] and appearance transformation [57], [58], and natural language texts mostly use back-translation approach [53] and spans technique [54]. Then, a model is pre-trained to identify whether the augmented samples are from the same original sample.

Recently, some studies [24], [27], [28], [59], [60] try to use contrastive learning approaches on unsupervised code representation learning. For example, Jain et al. [27] and Bui et al. [28] mainly use semantic-preserving program transformations to generate the functionally equivalent code snippets and pre-train the model to recognize semantically equivalent

and non-equivalent code snippets through contrastive learning techniques. These transformations including variable renaming (rename a variable with a random token), unused statement (insert a dead code snippet such as an unused declaration statement), permute statement (swap two statements which have no data dependency on each other), etc. Unlike the above-mentioned pre-trained technique, our model is based on multimodal contrastive learning with momentum encoders, which allow the model to learn the good representation based on samples in the current mini-batch and previous mini-batches. Furthermore, previous data augmentations require preserving the semantics of source code, whereas we use a simple yet effective dynamic masking technique that allows more flexible soft data augmentation.

## III. PROPOSED APPROACH

In this section, we illustrate our model CoCoSoDa for code search. The overall framework of CoCoSoDa is shown in Fig. 1. It is comprised of the following components:

- **Pre-trained code/query encoder** captures the semantic information of a code snippet or a natural language query and maps it into a high-dimensional embedding space.
- **Momentum code/query encoder** encodes the samples (code snippets or queries) of current and previous mini-batches to enrich the negative samples.
- **Soft data augmentation** is to dynamically mask or replace some tokens in a sample (code/query) to generate a similar sample as a form of data augmentation.
- **Multimodal contrastive learning** is used as the optimization objective and consists of inter-modal and intra-modal contrastive learning loss. They are used to minimize the distance of the representations of similar samples and maximize the distance of different samples in the embedding space.

We first illustrate our model with a concrete example shown in Fig. 1 and then introduce each component in details.

### A. An Illustrative Example

In this section, we introduce our model by an illustrative example shown in Fig. 1. On the left side is an example of soft data augmentation performing on the code snippet and query pair, and on the right side is the main architecture of our model. Specifically, *first*, at each iteration, we randomly perform one of four SoDa methods including dynamic masking (DM), dynamic replacement (DR), dynamic replacement of specified type (DRST), and dynamic masking of specified type (DMST), to generate the positive samples. DM and DR randomly sample 15% of tokens of a code snippet and replace each token with a [MASK] token or the type of the token. DRST and DMST sample all tokens of the specified type (such as operator, identifier) from a code snippet, and 15% of them are randomly replaced with the specified type or [MASK] token. For query, only DM is performed because other three SoDa methods require the type information of source code.

*Second*, we adopt the framework of momentum contrastive learning (MoCo) [30], and apply the multi-layer Transformer

encoder [61] as the backbone of the (momentum) encode because Transformer can effectively model and represent source code [22], [23]. We further extend MoCo for multi-modal representation learning by doubling the encoder and momentum encoder. Next, we use a large-scale code pre-trained model UniXcoder [24] to initialize the code/query encoder and momentum code/query encoder (Fig. 1), and then feed the original and augmented samples (codes or queries) to the encoders and momentum encoders, respectively, to obtain representations of samples. The momentum encoder can generate large and consistent representations of negative samples compared to the common encoder updated by back-propagation and the fixed encoder [30]. Therefore, at each iteration, CoCoSoDa can be trained to distinguish one sample from more other negative samples, so our approach can learn better representations of code snippets and queries.

*Third*, the multimodal contrastive learning consists of an intra-modal loss and an inter-modal loss is used to pull close similar representations and push apart dissimilar representations. Specifically, the intra-modal loss is used to learn an uniform distribution of representations of unimodal data by pulling in similar samples (codes or queries) and pushing away different samples. The inter-modal is use to learn the alignment of the multimodal data by pull close representations of the paired code snippet and query and push apart representations of the unpaired code snippet and query.

*Finally*, the well-trained model is used for code search. In detail, the code and query encoders map code snippets of codebase and the given query into a high-dimensional space, measure the similarity between them with cosine similarity, and return the most relevant code snippet based on the similarity.

### B. Pre-trained Encoder and Momentum Encoder

In this section, we introduce the base architecture, input samples, output representation and update mechanism of encoder and momentum encoder. The encoders and momentum encoders are all built on the multi-layer bidirectional Transformer Encoder [61]. As the pre-trained models such as UniXcoder [24] have achieved substantial improvement in code search, we initialize code and query encoder with parameters of UniXcoder. Following previous study [24], we average all the hidden states of the last layer as the whole sequence-level representation of query/code.

In the MoCo [30] framework, there is a momentum encoder encoding the samples of the current and previous mini-batches. Specifically, the momentum encoder maintains a queue by enqueueing the samples in the current mini-batch and dequeueing the samples in the oldest mini-batch. Here, we also use UniXcoder to initialize parameters of momentum code and query encoder. The difference of the update mechanism between the encoder and momentum encoder is that the encoder is updated by the back-propagation algorithm while the momentum encoder is updated by linear interpolation of the encoder and the momentum encoder. Thus, compared with the memory bank approach [62], which fixes and saves

the representations of all samples of the training dataset in advance, the momentum encoder can generate consistent representations and has been demonstrated to be effective [30]. For the end-to-end approach [19], [22], it has one encoder and takes other samples in the current mini-batch as negative samples. Thus, it requires a large mini-batch size in order to expand the number of negative samples [30]. For example, under the same computational resource such as A100-PCIE-80GB [63], up to 199 negative samples can be used in a mini-batch for end-to-end approaches, but our approach can use over 4,000 negative samples. Therefore, to support more negative samples, end-to-end approaches require larger memory computational resources than our approach.

We denote the parameters of the code encoder as  $\theta_{ce}$  and the momentum code encoder as  $\theta_{mce}$ , with parameters being the weights of UniXcoder. Therefore,  $\theta_{mce}$  is updated by:

$$\theta_{mce} = m\theta_{mce} + (1 - m)\theta_{ce} \quad (1)$$

where  $m \in [0, 1)$  is a momentum coefficient. Similarly, we denote the parameters of the query encoder and moment query encoder as  $\phi_{qe}$  and  $\phi_{mqe}$ . Then  $\phi_{mqe}$  is updated by:

$$\phi_{mqe} = m\phi_{mqe} + (1 - m)\phi_{qe} \quad (2)$$

Both  $\theta_{ce}$  and  $\phi_{qe}$  are learnable parameters and updated by the back-propagation algorithm.

### C. Soft Data Augmentation

In this section, we introduce soft data augmentation (SoDa) methods, which are simple data augmentation approaches without external constraints for source code or queries. We first introduce how to obtain soft data augmentation and then introduce how to use the augmented data.

The four SoDa methods are shown as follows.

- Dynamic Masking (DM): randomly sampling 15% of tokens of a code snippet and replace each token with a [MASK] token.
- Dynamic Replacement (DR): randomly sampling 15% of tokens of a code snippet and replace each token with the type of the token.
- Dynamic Replacement of Specified Type (DRST): sampling all tokens of a specified type (such as operator, identifier) from a code snippet, and 15% of them are randomly replaced with the specified type.
- Dynamic Masking of Specified Type (DMST): sampling all tokens of a specified type (such as operator, identifier) from a code snippet, and 15% of them are randomly masked.

Here, *dynamic* means that in data processing, the masking or replacement operation is performed at each iteration rather than only performed once [31]. It is worth noting that we randomly perform one of four SoDa methods for code snippets at each iteration. In addition, only DM is performed for queries because other three SoDa methods require the type information of source code.

We denote the SoDa module as  $G_{soda}$  which performs data transformation operations for the given input sequence to obtain the augmented data. Specifically, we first perform one of the four SoDa methods for the code snippets  $C = (c_1, \dots, c_{bs})$  and queries  $Q = (q_1, \dots, q_{bs})$  in a mini-batch by:

$$c_i^* = G_{soda}(c_i), \quad q_i^* = G_{soda}(q_i) \quad (i = 1, \dots, bs) \quad (3)$$

where  $c_i^*$  and  $q_i^*$  are the augmented samples of the code snippet  $c_i$  and query  $q_i$ , respectively.  $bs$  is mini-batch size. Then code snippet  $c_i$  and query  $q_i$  are fed into the code/query encoder and augmented samples  $c_k^*$  and  $q_k^*$  ( $k = 1, \dots, K$  and  $K$  is the queue size) in the current and previous mini-batches are fed to the momentum code/query encoder by:

$$\begin{aligned} \mathbf{v}_{c_i} &= f_{\theta_{ce}}(c_i), & \mathbf{v}_{c_k^*} &= f_{\theta_{mce}}(c_k^*) \\ \mathbf{v}_{q_i} &= f_{\theta_{qe}}(q_i), & \mathbf{v}_{q_k^*} &= f_{\theta_{mqe}}(q_k^*) \end{aligned} \quad (4)$$

where,  $\mathbf{v}_{c_i}$ ,  $\mathbf{v}_{q_i}$ ,  $\mathbf{v}_{c_k^*}$ , and  $\mathbf{v}_{q_k^*}$  are the final overall representations of the code snippet  $c_i$ , query  $q_i$ , augmented code snippet  $c_k^*$ , and augmented query  $q_k^*$ , respectively.

#### D. Multimodal Contrastive Learning

Multimodal contrastive learning consists of inter-modal and intra-modal loss function, and is used to optimize the parameters of the model. Specifically, given a query  $q_i$ , we denote the paired  $c_i$  or  $c_i^*$  as  $c_i^+$  and unpaired  $c_k^*$  as  $c_k^-$  ( $i = 1, \dots, bs$  and  $k = 1, \dots, K$ ). For the query  $q_i$ , with similarity measured by cosine similarity ( $sim(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$ ), we define the inter-modal and intra-modal contrastive learning loss [20], [64] as:

$$\begin{aligned} L_{q_i}^{inter} &= -\log \frac{e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{c_i^+})/\tau)}}{e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{c_i^+})/\tau)} + \sum_{k=1}^K e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{c_k^-})/\tau)}} \\ L_{q_i}^{intra} &= -\log \frac{e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{q_i^+})/\tau)}}{e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{q_i^+})/\tau)} + \sum_{k=1}^K e^{(sim(\mathbf{v}_{q_i}, \mathbf{v}_{q_k^-})/\tau)}} \end{aligned} \quad (5)$$

where  $\tau$  is the temperature hyperparameter [30], [62] and is set to 0.07 following previous works [27], [30]. Intuitively, the optimization objective of inter-modal loss function is to maximize the semantic similarity of the query and its paired code snippet and minimize the semantic similarity of the query and its unpaired code snippets. The intra-modal loss function is to learn the better representations of queries, where similar queries have closed representations and different queries have distinguishing representations. In the same way, for a code snippet  $c_i$ , we define the corresponding multimodal contrastive learning loss as:

$$\begin{aligned} L_{c_i}^{inter} &= -\log \frac{e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{q_i^+})/\tau)}}{e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{q_i^+})/\tau)} + \sum_{k=1}^K e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{q_k^-})/\tau)}} \\ L_{c_i}^{intra} &= -\log \frac{e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{c_i^+})/\tau)}}{e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{c_i^+})/\tau)} + \sum_{k=1}^K e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{c_k^-})/\tau)}} \end{aligned} \quad (6)$$

where  $q_i^+$  is the paired query of input code snippet  $c_i$ , and  $q_k^-$  denotes the unpaired query.

TABLE I  
DATASET STATISTICS.

Language	Training	Validation	Test	Candidate Code
Ruby	24,927	1,400	1,261	4,360
JavaScript	58,025	3,885	3,291	13,981
Java	164,923	5,183	10,955	40,347
Go	167,288	7,325	8,122	28,120
PHP	241,241	12,982	14,014	52,660
Python	251,820	13,914	14,918	43,827

The inter-modal and intra-modal loss function in a mini-batch can be obtained by:

$$L^{inter} = \sum_{i=1}^{bs} (L_{q_i}^{inter} + L_{c_i}^{inter}), \quad L^{intra} = \sum_{i=1}^{bs} (L_{q_i}^{intra} + L_{c_i}^{intra}) \quad (7)$$

To this end, the overall multimodal contrastive learning loss function for a mini-batch is:

$$L = \sum_{i=1}^{bs} (L^{inter} + L^{intra}) \quad (8)$$

We apply AdamW [65] to optimize the overall model.

#### E. Fine-tuning on Code Search

After being optimized by multimodal contrastive learning, the model can learn better representations of samples, where similar samples (code or queries) have similar representations and different samples have different representations. To further improve the performance of model on code search, following the most previous studies [22]–[26], [49], [59], we fine-tune it on the related training dataset by:

$$L^f = -\sum_{i=1}^{bs} \left[ \log \frac{e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{q_i})/\tau)}}{\sum_{j=1}^{bs} e^{(sim(\mathbf{v}_{c_i}, \mathbf{v}_{q_j})/\tau)}} \right] \quad (9)$$

where  $\mathbf{v}_{c_i}$  and  $\mathbf{v}_{q_i}$  are the overall semantic representations of the code snippet  $c_i$  and query  $q_i$ , respectively. They are obtained by code and query encoder, respectively.  $\tau$  is the temperature hyperparameter. Then we use the validation dataset to select the best model based on the MRR value (details in Sec. IV-D) and report scores on the test set in this paper.

## IV. EXPERIMENTAL DESIGN

### A. Datasets

We conduct experiments on a large-scale benchmark dataset CodeSearchNet [29] as used in Guo et al. [22]. It contains six programming languages, namely Ruby, JavaScript, Go, Python, Java, and PHP. This dataset is widely used in previous studies [11], [22]–[25], [29], [49], [66], [67]. The statistics of the dataset are shown in Table I. Following previous studies [10], [22], [68], the model is to retrieve the correct code snippets from the *Candidate Code* (the last column in Table I) for the given queries when performing the evaluation.

## B. Baselines

To evaluate the effectiveness of our approach, we compare CoCoSoDa with three IR-based methods [69]–[71], four deep end-to-end approaches [29] and ten pre-training-based approaches including three contrastive learning-related models.

- IR-based methods include **BOW** [69], **TF-IDF** [70] and **Jaccard** [71]. BOW and TF-IDF use bag-of-word and term frequency-inverse document frequency techniques, respectively, to extract the features from the input code snippets and queries and convert them into vectors. Then, they measure semantic similarities between code snippets and queries by cosine similarities between their corresponding vectors. Jaccard retrieves the similar code snippet for the given query according to the Jaccard similarity coefficient [71] between the code and query.
- **NBow**, **CNN**, **BiRNN** and **SelfAttn** [29] use various encoding models such as neural bag-of-words [72], 1D convolutional neural network [73], bi-directional GRU [74], and multi-head attention [61] to obtain the representations of code snippets and queries. And they measure the semantic similarity of representations using inner product.
- **RoBERTa** [31], **RoBERTa (code)** [23] are built on a multi-layer Transformer encoder [61] and pre-trained with MLM, which is to predict the masked tokens. The pre-trained datasets are natural language corpus [31] and source code corpus [29], respectively.
- **CodeBERT** [23] and **GraphCodeBERT** [22] are pre-trained on a large code corpus. The former is pre-trained with MLM and RTD, which uses a discriminator to identify replaced tokens. The latter considers the code structure information and is pre-trained tasks with MLM, data flow edge prediction, and node alignment.
- **Corder** [28] firstly uses a unimodal contrastive learning approach (only the code modality) to pre-train the model to recognize the semantically equivalent code snippets and then fine-tune it on the downstream tasks. As Corder does not release the implementation of semantic-preserving transformations and it is costly to implement these transformations for six programming languages, we only implemented for Java language because Java is the most studied language for code search [75]. To make a fair comparison, we use the unimodal contrastive learning technique to continually pre-train the UniXcoder and then fine-tune it on the code search task as the implementation of Corder.
- **ContraCode** [27] also applies contrastive learning to unsupervised code representation learning and conducts experiments on code summarization and type inference. Specifically, they first adopt some semantic-preserving program transformations to generate functional equivalence code snippets. Next, they pre-train a neural network model to identify functionally equivalent code snippets among many distractors. Finally, they fine-tune the pre-trained model to perform downstream tasks. We use the pre-trained ContraCode as the code/query encoder and optimize it using Eq. 9 for code search.

- **CodeT5** [25], **PLBART** [26], **SPT-Code** [49] are sequence-to-sequence code pre-trained models. The first is pre-trained with three identifier-aware pre-training tasks to enable the model to identify identifiers in source code or recover masked identifiers. The second is pre-trained with denoising autoencoding, which is to reconstruct the corrupted input code sequence. The third takes source code, corresponding AST and paired summarization as input and is pre-trained with three code-specific tasks.
- **SyncoBERT** [59] and **UniXcoder** [24] are multi-modal contrastive pre-training for code representation. SyncoBERT takes source code, AST and summarization as input and pre-trained with identifier prediction and AST edge prediction to learn the lexical and syntactic knowledge of source code. UniXcoder takes two-modality data, the summarization and simplified AST of source code, as input and is pre-trained with MLM, unidirectional language modeling, denoising autoencoder, and two contrastive learning-related tasks.

More details about baselines can be found in Appendix of the replication package [76]. In our experiments, we train the four deep end-to-end approaches from scratch, and for the ten pre-trained approaches, we initialize them with the pre-trained models and fine-tune (or continually pre-train and fine-tune) them according to the descriptions in their original papers or their released source code.

## C. Experimental Settings

Following UniXcoder [24], we use Transformer with 12 layers, 768 dimensional hidden states, and 12 attention heads. The vocabulary sizes of code and queries are set to 51,451. Max sequence lengths of code snippets and queries are 128 and 256, respectively. For optimizer, we use AdamW with the learning rate  $2e-5$ . Following previous studies [22], [23], [68], the code encoder and query encoder share parameters to reduce the number of total parameters. Following MoCo [30], the temperature hyperparameter  $\tau$  is set as 0.07 and momentum coefficient  $m$  is 0.999. The queue size and batch size are set to 4096 and 128, respectively. The training step of multimodal contrastive learning stage is 100K and the maximum epochs of fine-tune stage is 5. In addition, we run the experiments 3 times with random seeds 0,1,2 and display the mean value in the paper. All experiments are conducted on a machine with 220 GB main memory and Tesla A100 80GB GPU.

## D. Evaluation Metrics

We measure the performance of our approach using four metrics: mean reciprocal rank (MRR) and top-k recall ( $R@k$ ,  $k=1,5,10$ ), which are widely used in previous studies [10], [11], [13]–[20], [20], [20]–[23], [29]. **MRR** is the average of reciprocal ranks of the correct code snippets for given queries  $Q$ .  **$R@k$**  measures the percentage of queries that the paired code snippets exist in the top-k returned ranked lists. They are calculated as follows:

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{Rank_i}, \quad R@k = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \delta(Rank_i \leq k) \quad (10)$$

TABLE II  
PERFORMANCE OF DIFFERENT APPROACHES. JS IS SHORT FOR  
JAVASCRIPT. STATISTICAL SIGNIFICANCE OF EXPERIMENTS:  $p < 0.01$ .

	Model	Ruby	JS	Go	Python	Java	PHP	Avg.
IR-Based models	BOW	0.230	0.184	0.350	0.222	0.245	0.193	0.237
	TF-IDF	0.239	0.204	0.363	0.240	0.262	0.215	0.254
	Jaccard	0.220	0.191	0.345	0.243	0.235	0.182	0.236
Deep end-to-end models	NBow	0.162	0.157	0.330	0.161	0.171	0.152	0.189
	CNN	0.276	0.224	0.680	0.242	0.263	0.260	0.324
	BiRNN	0.213	0.193	0.688	0.290	0.304	0.338	0.338
	SelfAtt	0.275	0.287	0.723	0.398	0.404	0.426	0.419
Pre-trained models	RoBERTa	0.587	0.523	0.855	0.590	0.605	0.561	0.620
	RoBERTa (code)	0.631	0.57	0.864	0.621	0.636	0.581	0.650
	CodeBERT	0.679	0.621	0.885	0.672	0.677	0.626	0.693
	GraphCodeBERT	0.703	0.644	0.897	0.692	0.691	0.649	0.713
	Corder	-	-	-	-	0.727	-	-
	ContraCode	-	0.688	-	-	-	-	-
	PLBART	0.675	0.616	0.887	0.663	0.663	0.611	0.685
	CodeT5	0.719	0.655	0.888	0.698	0.686	0.645	0.715
	SyncoBERT	0.722	0.677	0.913	0.724	0.723	0.678	0.740
	UniXcoder	0.740	0.684	0.915	0.720	0.726	0.676	0.744
SPT-Code	0.701	0.641	0.895	0.699	0.700	0.651	0.715	
Our	CoCoSoDa	<b>0.818</b> ↑10.54%	<b>0.764</b> ↑11.7%	<b>0.921</b> ↑0.66%	<b>0.757</b> ↑5.14%	<b>0.763</b> ↑5.10%	<b>0.703</b> ↑3.99%	<b>0.788</b> ↑5.91%

where  $Rank_i$  is the rank of the paired code snippet related to the  $i$ -th query.  $\delta$  is an indicator function that returns 1 if  $Rank_i \leq k$  otherwise returns 0.

## V. EXPERIMENTAL RESULTS

### A. RQ1: What Is the Effectiveness of CoCoSoDa?

1) *Overall results:* We evaluate the effectiveness of our model CoCoSoDa by comparing it to three IR-based models, four recent deep end-to-end code search models and ten pre-trained models introduced in Sec. IV-B on the CodeSearchNet dataset with six programming languages. As there is no released source code for Corder and it is costly to reproduce this approach for six programming languages, we reproduce Corder on Java language for comparison. The experimental results are shown in Table II. We present the results under MRR metric only due to space limitation. We put results under other metrics in Appendix of the replication package [76]. Conclusions that hold on MRR also hold for other metrics.

We can see that deep end-to-end models outperform IR-based models because they can learn better semantic relations between codes and queries [10], [14], [20]. The ten pre-trained models (the third row of Table II) perform better than the four deep end-to-end models (the second row of Table II) trained from scratch, which shows the effectiveness of the pre-training technique. Since GraphCodeBERT considers the data flow information of source code, it performs better than RoBERTa, RoBERTa (code) and CodeBERT. UniXcoder and SyncoBERT consider the structure information of source code and perform better than other approaches on average. Corder and ContraCode perform best among baselines because they pre-train the model to recognize the functionally equivalent code snippets from many distractors and can learn a better code representation. CoCoSoDa takes UniXcoder as the base code/query encoder, continually optimizes it with multimodal

Listing 1. The top-1 result returned by CoCoSoDa.

```
public static byte[] hexStringToByte(String hexString) {
    try {
        return Hex.decodeHex(hexString.toCharArray());
    } catch (DecoderException e) {
        throw new UnexpectedException(e);
    }
}
```

Listing 2. The top-1 result returned by GraphCodeBERT.

```
public static String toHexString(final byte[] bytes) {
    char[] chars = new char[bytes.length * 2];
    int i = 0;
    for (byte b : bytes) {
        chars[i++] = CharUtil.int2hex((b & 0xF0) >> 4);
        chars[i++] = CharUtil.int2hex(b & 0x0F);
    }
    return new String(chars);
}
```

Fig. 2. The top-1 code returned by CoCoSoDa and GraphCodeBERT for the query “Transform a hexadecimal String to a byte array.” on Java language.

Listing 3. The top-1 result returned by CoCoSoDa.

```
def get_items(self):
    reader = csv.reader(self.source)
    headers = reader.next()
    for row in reader:
        if not row:
            continue
        yield dict(zip(headers, row))
```

Listing 4. The top-1 result returned by UniXcoder.

```
def iterrows(lines_or_file, namedtuples=False, dicts=False,
            encoding='utf-8', **kw):
    if namedtuples and dicts:
        raise ValueError('either namedtuples or dicts can
            be chosen as output format')
    elif namedtuples:
        _reader = NamedTupleReader
    elif dicts:
        _reader = UnicodeDictReader
    else:
        _reader = UnicodeReader
    with _reader(lines_or_file, encoding=encoding, **fix_kw
                (kw)) as r:
        for item in r:
            yield item
```

Fig. 3. The top-1 code returned by CoCoSoDa and UniXcoder for the query “Iterator to read the rows of the CSV file.” on Python language.

contrastive learning and soft data augmentation and performs best among all approaches.

2) *Case study:* Next, we show some cases to demonstrate the effectiveness of our model CoCoSoDa. For each case, we only show the result of our approach and the best baseline, which is GraphCodeBERT for the 1st case and UniXcoder for the 2nd case.

Fig. 2 shows the results returned by CoCoSoDa and GraphCodeBERT for the query “Transform a hexadecimal String to a byte array.” from the Java dataset. The query includes two operation objects: “hexadecimal String” and “byte array” and one action “Transform”. To implement the functionality of the query, we usually take the “hexadecimal String” as an input parameter and use “toXXX(…)” to perform the “Transform” action. Our model CoCoSoDa can successfully understand the semantics of the whole query and code snippet and return the correct result, while GraphCodeBERT cannot. This is because the code representation obtained by GraphCodeBERT is affected by token-level semantics such as `String` and `byte`, thereby, returning the code snippet which has many similar tokens with the query but a completely opposite semantic: “Converts bytes to hex string.”

In Fig. 3, we compare the results returned by CoCoSoDa and UniXcoder for the query “Iterator to read the rows of

TABLE III  
ABLATION STUDY OF CoCoSoDa ON MRR.

Model	Ruby	JS	Go	Python	Java	PHP	Avg.
<b>CoCoSoDa</b>	<b>0.818</b>	<b>0.764</b>	<b>0.921</b>	<b>0.757</b>	<b>0.763</b>	<b>0.703</b>	<b>0.788</b>
w/o DR	0.794	0.714	0.905	0.730	0.743	0.688	0.762
w/o DM	0.801	0.754	0.906	0.744	0.756	0.69	0.775
w/o DRST	0.797	0.714	0.905	0.730	0.743	0.685	0.762
w/o DMST	0.788	0.715	0.902	0.731	0.746	0.686	0.761
w/o all SoDas	0.775	0.711	0.907	0.736	0.738	0.683	0.758
w/o inter-modal loss	0.776	0.703	0.906	0.729	0.739	0.674	0.755
w/o intra-modal loss	0.780	0.705	0.903	0.727	0.744	0.680	0.756

the CSV file.” in the Python dataset. CoCoSoDa returns the correct code snippet, while UniXcoder returns the code snippet with another semantics “Yield a generator over the rows.”. Our model can accurately understand the intent of the query and return the relevant code snippet, which first reads a CSV file and yields an iterator to read the rows of it. UniXcoder returns a partially relevant code snippet which only yields an iterator to read the rows of files, missing reading a CSV file.

**Summary.** Our approach significantly outperforms baselines on six programming languages in terms of four metrics. Case studies further demonstrate the advantages of CoCoSoDa in code search.

### B. RQ2: How Much Do Different Components Contribute?

In this section, we study the contribution of each component of our approach CoCoSoDa. It includes multimodal contrastive learning such as intra-modal loss and inter-modal loss and four SoDa approaches (DR, DM, DRST and DMST) introduced in Sec. III-C. Specifically, we remove one component (such as DM) of CoCoSoDa each time and then study the performance of the ablated model. The experimental results are shown in Table III. “w/o one component” means to remove this component. For example, CoCoSoDa w/o DR and inter-modal loss means to drop the SoDa method DR and inter-modal loss function  $L^{intra}$  (Eq. 7), respectively.

From the Table III, we can see that the performance of the model drops after removing any one component. It demonstrates that each component plays an important role in the code search model. Especially, the performance of CoCoSoDa w/o all SoDas, intra-modal and inter-modal contrastive learning drops obviously. This is because data augmentation can increase data diversity. The inter-modal contrastive learning, which aims to learn the alignment of code snippets and queries, can pull together the paired code and query and push apart the unpaired code and query. Intra-modal contrastive learning aims to learn a uniform distribution of representation of unimodal samples (code snippets or queries) and can improve the generalization performance of code search model.

**Summary.** The ablation study shows the effectiveness of multimodal contrastive learning including intra-modal loss and inter-modal loss and four SoDa approaches including DR, DM, DRST and DMST.

### C. RQ3: What Is the Performance of Our Approach on Other Pre-trained Models?

We further study the performance of our approach on other three pre-trained models introduced in Sec. IV-B, including a natural language pre-trained model RoBERTa and two source code pre-trained models CodeBERT and GraphCodeBERT. Specifically, we use these pre-trained models as the code/query encoders and momentum code/query encoders in Fig. 1. For the input code snippet and query sequence, we average hidden states of the last layer as the overall representations of the code snippet or query. The similarities of representations are measured by the cosine similarity. Other experimental settings are same as in Sec. IV-C.

The results are shown in Table IV. CoCoSoDa<sub>RoBERTa</sub> means using RoBERTa as the code/query encoders and momentum code/query encoders in our framework. Overall, we can see that CoCoSoDa<sub>RoBERTa</sub>, CoCoSoDa<sub>CodeBERT</sub> and CoCoSoDa<sub>GraphCodeBERT</sub> obviously outperform RoBERTa, CodeBERT and GraphCodeBERT, respectively on all six programming languages in terms of four metrics. These results demonstrate that our approach can be generalized to other pre-trained models and boost their performance. Besides, CoCoSoDa<sub>RoBERTa</sub>, which is pre-trained on the natural language corpus and fine-tuned with our method, achieves comparable performance with CodeBERT on Go dataset. CoCoSoDa<sub>GraphCodeBERT</sub> achieves comparable performance with UniXcoder on JavaScript, Java and Python. CoCoSoDa<sub>GraphCodeBERT</sub> also slightly outperforms UniXcoder on Ruby dataset.

**Summary.** Our approach is orthogonal to the pre-trained technique on the performance improvement for code search tasks and can obviously boost the performance of existing pre-trained models.

### D. RQ4: What Is the Impact of Different Hyperparameters?

In this section, we study the impact of different hyperparameters: learning rate, momentum coefficient  $m$ , masking ratio  $r$ , and temperature hyperparameter  $\tau$ . We study different hyperparameters in the typical range, which covers all experimental settings of previous studies [22]–[25], [28], [31], [49], [59] and the experimental results are shown in Fig. 4. From the results of varying learning rate (the top left of Fig. 4), we can see that performance is generally stable for small learning rate [77] (from  $5e^{-6}$  to  $7e^{-5}$ ). The learning rates that are larger than  $7e^{-5}$  have obvious impacts on the model performance. The results of different momentum coefficient  $m$  are shown in the top right of Fig. 4. We can see that performance increases when the momentum coefficient  $m$  becomes larger. This is because a large momentum coefficient is beneficial to obtain the consistent representation for the queue [30]. The momentum coefficient that is smaller than 0.910 has a significant impact on performance. These findings are consistent with the previous work [30]. From the results of varying masked ratio  $r$  (the bottom left of Fig. 4), we can see that the performance is insensitive to the masked ratio  $r$  when the masked ratio  $r$  is between 5% and 20%. A larger

TABLE IV  
RESULTS ON OTHER PRE-TRAINED MODELS. CoCoSoDa<sub>GRAPH</sub> IS SHORT FOR CoCoSoDa<sub>GRAPH</sub>CodeBERT. THE IMPROVED PERCENTAGES ARE SHOWN IN PARENTHESES.

PL	Metric	RoBERTa	CoCoSoDa <sub>RoBERTa</sub>	CodeBERT	CoCoSoDa <sub>CodeBERT</sub>	GraphCodeBERT	CoCoSoDa <sub>Graph</sub>
Ruby	MRR	0.587	<b>0.640</b> (↑9.03%)	0.679	<b>0.723</b> (↑6.48%)	0.703	<b>0.752</b> (↑6.97%)
	R@1	0.469	<b>0.533</b> (↑13.65%)	0.583	<b>0.618</b> (↑6.00%)	0.607	<b>0.655</b> (↑7.91%)
	R@5	0.717	<b>0.764</b> (↑6.56%)	0.800	<b>0.852</b> (↑6.50%)	0.824	<b>0.875</b> (↑6.19%)
	R@10	0.785	<b>0.825</b> (↑5.10%)	0.853	<b>0.904</b> (↑5.98%)	0.872	<b>0.916</b> (↑5.05%)
JavaScript	MRR	0.523	<b>0.559</b> (↑6.88%)	0.621	<b>0.648</b> (↑4.35%)	0.644	<b>0.682</b> (↑5.90%)
	R@1	0.413	<b>0.460</b> (↑11.38%)	0.514	<b>0.545</b> (↑6.03%)	0.538	<b>0.582</b> (↑8.18%)
	R@5	0.652	<b>0.673</b> (↑3.22%)	0.752	<b>0.772</b> (↑2.66%)	0.774	<b>0.806</b> (↑4.13%)
	R@10	0.730	<b>0.744</b> (↑1.92%)	0.814	<b>0.839</b> (↑3.07%)	0.834	<b>0.866</b> (↑3.84%)
Go	MRR	0.855	<b>0.881</b> (↑3.04%)	0.885	<b>0.905</b> (↑2.26%)	0.897	<b>0.907</b> (↑1.11%)
	R@1	0.800	<b>0.829</b> (↑3.62%)	0.837	<b>0.859</b> (↑2.63%)	0.858	<b>0.861</b> (↑0.35%)
	R@5	0.926	<b>0.945</b> (↑2.05%)	0.944	<b>0.962</b> (↑1.91%)	0.954	<b>0.962</b> (↑0.84%)
	R@10	0.949	<b>0.965</b> (↑1.69%)	0.962	<b>0.975</b> (↑1.35%)	0.972	<b>0.978</b> (↑0.62%)
Python	MRR	0.590	<b>0.629</b> (↑6.61%)	0.672	<b>0.690</b> (↑2.68%)	0.692	<b>0.714</b> (↑3.18%)
	R@1	0.480	<b>0.523</b> (↑8.96%)	0.574	<b>0.589</b> (↑2.61%)	0.594	<b>0.614</b> (↑3.37%)
	R@5	0.727	<b>0.756</b> (↑3.99%)	0.792	<b>0.809</b> (↑2.15%)	0.813	<b>0.834</b> (↑2.58%)
	R@10	0.793	<b>0.818</b> (↑3.15%)	0.850	<b>0.867</b> (↑2.00%)	0.866	<b>0.888</b> (↑2.54%)
Java	MRR	0.605	<b>0.635</b> (↑4.96%)	0.677	<b>0.705</b> (↑4.14%)	0.691	<b>0.721</b> (↑4.34%)
	R@1	0.499	<b>0.531</b> (↑6.41%)	0.580	<b>0.606</b> (↑4.48%)	0.592	<b>0.624</b> (↑5.41%)
	R@5	0.737	<b>0.762</b> (↑3.39%)	0.796	<b>0.826</b> (↑3.77%)	0.817	<b>0.843</b> (↑3.18%)
	R@10	0.796	<b>0.820</b> (↑3.02%)	0.852	<b>0.878</b> (↑3.05%)	0.865	<b>0.890</b> (↑2.89%)
PHP	MRR	0.561	<b>0.598</b> (↑6.60%)	0.626	<b>0.647</b> (↑3.35%)	0.649	<b>0.668</b> (↑2.93%)
	R@1	0.450	<b>0.490</b> (↑8.89%)	0.520	<b>0.538</b> (↑3.46%)	0.545	<b>0.561</b> (↑2.94%)
	R@5	0.694	<b>0.730</b> (↑5.19%)	0.753	<b>0.779</b> (↑3.45%)	0.785	<b>0.798</b> (↑1.66%)
	R@10	0.764	<b>0.797</b> (↑4.32%)	0.814	<b>0.844</b> (↑3.69%)	0.832	<b>0.863</b> (↑3.73%)

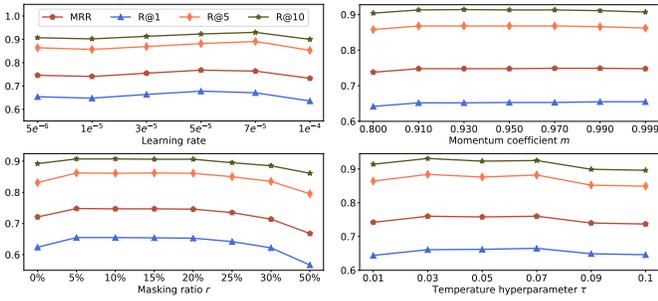


Fig. 4. The impact of different hyperparameters.

masked ratio such as 50% brings considerable performance degradation. It is reasonable because the larger masked ratio causes the code snippet to lose too much information. The results of different temperature hyperparameter  $\tau$  are shown in the bottom right of Fig. 4. We can see that performance is stable when the temperature hyperparameter  $\tau$  varies from 0.03 to 0.07.

**Summary.** In general, our model is stable over a range of hyperparameter values (learning rate is from  $5e^{-6}$  to  $7e^{-5}$ , momentum coefficient is between 0.910 and 0.999, masked ratio  $r$  is from 5% to 20%, and temperature hyperparameter  $\tau$  varies from 0.03 to 0.07).

#### E. RQ5: Why Does Our Model CoCoSoDa Work?

The advantages of CoCoSoDa mainly come from soft data augmentation and multimodal momentum contrastive learning. The soft data augmentation transforms the input sequence with the masking or replacement mechanism and generates a

similar sample as the “positive” example. It can help the model learn a representation of the code snippet and query from a global (sequence-level) view rather than simply aggregate the token-level semantic. Thus, our model tends to return code snippets related to the sequence-level functionality rather than the token-level similarities. The momentum mechanism enlarges negative samples, allowing to distinguish one sample from more negative samples at each iteration. For multimodal contrastive learning, the inter-modal contrastive learning can pull together the representations of the code-query pair and push apart the representations of queries and many unpaired code snippets, while intra-modal contrastive learning can learn a uniform distribution of representations in terms of unimodal data (code snippets or queries). Therefore, CoCoSoDa can learn better representations of code and queries and perform well on code search. We further explore why CoCoSoDa works through quantitative and qualitative analysis.

1) *Quantitative analysis:* We explore to understand reasons behind good performance of our approach through  $\ell_{\text{align}}$  and  $\ell_{\text{uniform}}$  [78], which are usually used as indicators to reflect the quality of representation learned by contrastive learning techniques [52], [78]–[80]. Mathematically, they are defined as:

$$\begin{aligned} \ell_{\text{align}} &= \mathbb{E}_{(x,y) \sim D_{\text{paired}}} \|f(x) - f(y)\|_2^2 \\ \ell_{\text{uniform}} &= \log \mathbb{E}_{(x,y) \stackrel{\text{i.i.d.}}{\sim} D} [e^{-2\|f(x) - f(y)\|_2^2}] \end{aligned} \quad (11)$$

where  $(x, y) \sim D_{\text{paired}}$  means the  $x$  and  $y$  are paired samples, while  $(x, y) \stackrel{\text{i.i.d.}}{\sim} D$  mean that  $x$  and  $y$  are independent identically distributed.  $f(x)$  and  $f(y)$  are learned representations and  $\|f(x) - f(y)\|_2^2$  represent the 2-norm of the distance between them. From the Eq. 11, we can know that  $\ell_{\text{align}}$  always is a non-negative number, while  $\ell_{\text{uniform}}$  is a non-

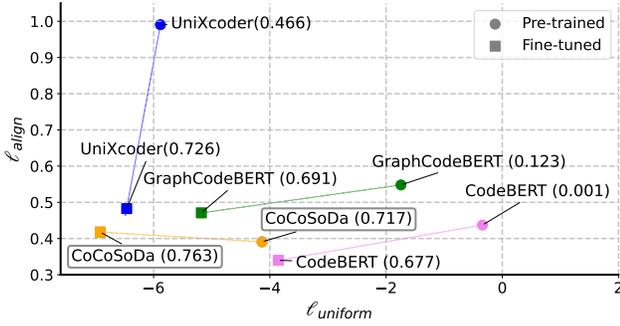


Fig. 5.  $\ell_{\text{align}}-\ell_{\text{uniform}}$  plot of different models. Circles and squares represent pre-trained and fine-tuned models, respectively. Identical models are marked with the same color. MRR scores are shown in parentheses.

positive number. In terms of our task,  $\ell_{\text{align}}$  reflects the degree of alignment of representations of code-query pairs. The closer distances of the paired code snippets and queries are, the smaller the value (close to zero) of  $\ell_{\text{align}}$  is.  $\ell_{\text{uniform}}$  reflects the uniformity of the distribution of representations of all code snippets or queries. The more uniform the distribution is, the smaller the value (close to negative infinity) of  $\ell_{\text{align}}$  is. In the extreme case where all representations of code snippets and queries are the same, both values of  $\ell_{\text{align}}$  and  $\ell_{\text{uniform}}$  are zero. That is to say, learned representations have perfect alignment but extremely poor uniformity. In fact, a model which can learn representations with both better alignment (lower  $\ell_{\text{align}}$ ) and uniformity (lower  $\ell_{\text{uniform}}$ ) can generally achieve better performance [52], [78].

We show  $\ell_{\text{align}}-\ell_{\text{uniform}}$  plot in Fig. 5. We can find that (1) pre-trained CodeBERT and GraphCodeBERT have better alignment but poor uniformity, and perform not well. This is because representations learned by them are very similar and cannot generalize well. UniXcoder has better uniformity and performs better than them. Pre-trained CoCoSoDa has both better alignment and uniformity than CodeBERT and GraphCodeBERT and performs best among four pre-trained models. (2) After being fine-tuned, UniXcoder, CodeBERT and GraphCodeBERT all have lower  $\ell_{\text{uniform}}$  and  $\ell_{\text{align}}$  and significantly outperform pre-trained ones individually. Our fine-tuned CoCoSoDa generally preserves the alignment and has a better uniform. Therefore, fine-tuned CoCoSoDa further improves the performance of pre-trained it.

2) *Qualitative analysis*: We also visualize learned representations to help intuitively understand why CoCoSoDa works well. Specifically, first, we randomly sample  $X$  ( $X=100, 200, 300, \text{ or } 400$ ) code-query pairs from Java dataset with ten different random seeds from 0 to 9. We show the result with  $X=300$  with a random seed of 3 in the paper, and other visualized results are put in Appendix of replication package [76] due to space limitation. The following findings and conclusions hold for different  $X$  and random seeds. Second, we feed sampled pairs including code snippets and queries to well-trained CoCoSoDa and UniXcoder individually and obtain their representations. Third, we apply T-SNE [81] to

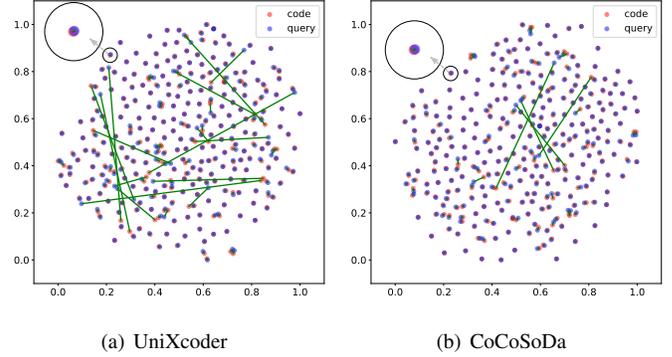


Fig. 6. T-SNE visualization of representations of code snippets and queries. (a) and (b) are the representations learned by UniXcoder and CoCoSoDa, respectively. Code is in orange and query is in blue. The distance of paired code and query is indicated by a green line.

TABLE V  
THE PERFORMANCE OF DIFFERENT APPROACHES UNDER THE ZERO-SHORT EXPERIMENTAL SETTING EVALUATED BY MRR SCORES.

Model	Ruby	JS	Go	Python	Java	PHP	Avg.
CodeBERT	0.002	0.001	0.002	0.001	0.001	0.001	0.001
CodeT5	0.006	0.003	0.009	0.001	0.001	0.001	0.004
GraphCodeBERT	0.238	0.111	0.209	0.137	0.123	0.120	0.156
UniXcoder	0.576	0.442	0.648	0.447	0.466	0.373	0.492
CoCoSoDa	<b>0.786</b>	<b>0.709</b>	<b>0.881</b>	<b>0.696</b>	<b>0.717</b>	<b>0.640</b>	<b>0.738</b>
	↑36.5%	↑60.4%	↑36.0%	↑55.7%	↑53.9%	↑71.6%	↑50.0%

reduce the dimensionality of obtained representations into 2D and visualize dimensionality-reduced representations in Fig. 6. In detail, Fig. 6(a) and Fig. 6(b) show representations learned by UniXcoder and CoCoSoDa, respectively. The code snippet and query are marked in orange and blue, respectively. The distance between the paired code and query is indicated by a green line. From the Fig. 6, we can see that (1) there are many very short green lines in Fig. 6 (a) and (b), which means that both CoCoSoDa and UniXcoder can map most paired code snippets and queries into close embeddings. (2) Fig. 6(a) have more long green lines than Fig. 6(b). It indicates that UniXcoder has more cases than CoCoSoDa, where pairwise representations of queries and code snippets are far away from each other. In the future, we will conduct error analysis to study the paired samples with long green lines. In summary, the visualization results intuitively show that CoCoSoDa can learn better representations than UniXcoder.

**Summary.** Our model can learn a more uniform distribution of the representations of unimodal data (code snippets or queries), and the learned representations of paired code snippets and queries can be well aligned.

## VI. DISCUSSION

### A. The Performance of CoCoSoDa Without Being Fine-Tuned

To further study the effectiveness of CoCoSoDa, we evaluate different pre-trained models under the zero-short experimental setting, where models are directly used for evaluation

without fine-tuning them. The experimental results are shown in Table V. We can see that pre-trained CodeBERT, CodeT5, and GraphCodeBERT without fine-tuning performs poorly due to the representation degeneration problem [82], [83]. That is, the high-frequent tokens dominate the sequence representation [82], resulting in the poor sequence-level semantic representation of the code snippet and query. UniXcoder and our models adopt the contrastive learning related technique to obtain better representations of code snippets and queries, and perform better than the other two models. Especially, our approach can help the model to learn a uniform distribution of representation of unimodal data, and learn a good alignment for multimodal data. Therefore, our model outperforms other pre-trained models by about 50% on average MRR scores even more than 70% on PHP language. Furthermore, the performance of CoCoSoDa under zero-shot setting even exceeds many fine-tuned models such as CodeBERT, CodeT5 and GraphCodeBERT (Table II) on average MRR. In summary, our pre-trained model performs better than other pre-trained models when all of them are not fine-tuned. Furthermore, the performance of CoCoSoDa without being fine-tuned is even better than many fine-tuned models on average MRR.

### B. Limitations & Threats to Validity

Although CoCoSoDa has an overall advantage, our model could still return inaccurate results, especially for the code snippets that use the third-library API or self-defined methods. This is because CoCoSoDa only considers the information of the code snippet itself rather than other contexts such as other methods in the enclosing class or project [66], [84]. In our future work, more contextual information (such as enclosing class/project and called API/methods) could be considered in our model to further improve the performance of CoCoSoDa.

We also identify the following threats to our approach:

*Programming Languages.* Due to the heavy effort to evaluate the model on all programming languages, we conduct our experiment with as many programming languages as possible on the existing build datasets. Our model on different programming languages would have different results. In the future, we will evaluate the effectiveness of our approach with more other programming languages.

*Pre-trained models.* To demonstrate that our approach is orthogonal to the pre-trained technique on the performance improvements for code search, we have adopted and evaluated our approach on four pre-trained models including a natural language pre-trained model RoBERTa and three source code pre-trained models CodeBERT, GraphCodeBERT and UniXcoder. It remains to be verified whether or not the proposed approach is applicable to other pre-trained models such as GPT [85] and T5 [25].

*Evaluated benchmark.* The paired code snippet is usually used as the correct result for the given query. In fact, some unpaired code snippets also answer the given query. In the future, we will invite some developers to manually score the semantical correlation between the arbitrary code snippet and query and build a high-quality code search benchmark.

## VII. CONCLUSION

In this paper, we present CoCoSoDa, which leverages multimodal momentum contrastive learning and soft data augmentation for code search. It can help the model learn effective representations by pulling together representations of code-query pairs and pushing apart the unpaired code snippets and queries. We conduct extensive experiments on a large-scale benchmark dataset with six programming languages and the results confirm its superiority. In our future work, more contextual information (such as enclosing class/project) could be considered in our model to further improve the performance of CoCoSoDa. Replication package including datasets, source code, and Appendix is available at <https://github.com/DeepSoftwareAnalytics/CoCoSoDa>.

## ACKNOWLEDGEMENT

We thank reviewers for their valuable comments on this work. This research was supported by National Key R&D Program of China (No. 2017YFA0700800) and Fundamental Research Funds for the Central Universities under Grant xtr072022001. We would like to thank Jiaqi Guo for their valuable suggestions and feedback during the work discussion process.

## REFERENCES

- [1] J. Singer, T. C. Lethbridge, N. G. Vinson, and N. Anquetil, "An examination of software engineering work practices," in *CASCON*. IBM, 1997, p. 21.
- [2] L. Nie, H. Jiang, Z. Ren, Z. Sun, and X. Li, "Query expansion based on crowd knowledge for code search," *IEEE Trans. Serv. Comput.*, vol. 9, no. 5, pp. 771–783, 2016.
- [3] C. McMillan, M. Grechanik, D. Poshyanyk, Q. Xie, and C. Fu, "Portfolio: finding relevant functions and their usage," in *ICSE*. ACM, 2011, pp. 111–120.
- [4] S. P. Reiss, "Semantics-based code search," in *ICSE*. IEEE, 2009, pp. 243–253.
- [5] F. Zhang, H. Niu, I. Keivanloo, and Y. Zou, "Expanding queries for code search using semantically related API class-names," *IEEE Trans. Software Eng.*, vol. 44, no. 11, pp. 1070–1082, 2018.
- [6] M. Allamanis, E. T. Barr, P. T. Devanbu, and C. Sutton, "A survey of machine learning for big code and naturalness," *ACM Comput. Surv.*, vol. 51, no. 4, pp. 81:1–81:37, 2018.
- [7] E. Linstead, S. K. Bajracharya, T. C. Ngo, P. Rigor, C. V. Lopes, and P. Baldi, "Sourcerer: mining and searching internet-scale software repositories," *Data Min. Knowl. Discov.*, vol. 18, no. 2, pp. 300–336, 2009.
- [8] M. Lu, X. Sun, S. Wang, D. Lo, and Y. Duan, "Query expansion via wordnet for effective code search," in *SANER*. IEEE Computer Society, 2015, pp. 545–549.
- [9] F. Lv, H. Zhang, J. Lou, S. Wang, D. Zhang, and J. Zhao, "Codehow: Effective code search based on API understanding and extended boolean model (E)," in *ASE*. IEEE Computer Society, 2015, pp. 260–270.
- [10] X. Gu, H. Zhang, and S. Kim, "Deep code search," in *ICSE*. ACM, 2018, pp. 933–944.
- [11] L. Du, X. Shi, Y. Wang, E. Shi, S. Han, and D. Zhang, "Is a single model enough? mucos: A multi-model ensemble learning approach for semantic code search," in *CIKM*. ACM, 2021, pp. 2994–2998.
- [12] J. Cambroner, H. Li, S. Kim, K. Sen, and S. Chandra, "When deep learning met code search," in *Proceedings of the ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/SIGSOFT FSE 2019, Tallinn, Estonia, August 26-30, 2019*, M. Dumas, D. Pfahl, S. Apel, and A. Russo, Eds. ACM, 2019, pp. 964–974. [Online]. Available: <https://doi.org/10.1145/3338906.3340458>

- [13] X. Ling, L. Wu, S. Wang, G. Pan, T. Ma, F. Xu, A. X. Liu, C. Wu, and S. Ji, "Deep graph matching and searching for semantic code retrieval," *ACM Trans. Knowl. Discov. Data*, vol. 15, no. 5, pp. 88:1–88:21, 2021. [Online]. Available: <https://doi.org/10.1145/3447571>
- [14] W. Li, H. Qin, S. Yan, B. Shen, and Y. Chen, "Learning code-query interaction for enhancing code searches," in *ICSME*. IEEE, 2020, pp. 115–126.
- [15] Q. Zhu, Z. Sun, X. Liang, Y. Xiong, and L. Zhang, "Ocor: An overlapping-aware code retriever," in *ASE*, 2020.
- [16] J. Shuai, L. Xu, C. Liu, M. Yan, X. Xia, and Y. Lei, "Improving code search with co-attentive representation learning," in *ICPC '20: 28th International Conference on Program Comprehension, Seoul, Republic of Korea, July 13-15, 2020*. ACM, 2020, pp. 196–207. [Online]. Available: <https://doi.org/10.1145/3387904.3389269>
- [17] W. Ye, R. Xie, J. Zhang, T. Hu, X. Wang, and S. Zhang, "Leveraging code generation to improve code retrieval and summarization via dual learning," in *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, Y. Huang, I. King, T. Liu, and M. van Steen, Eds. ACM / IW3C2, 2020, pp. 2309–2319. [Online]. Available: <https://doi.org/10.1145/3366423.3380295>
- [18] R. Haldar, L. Wu, J. Xiong, and J. Hockenmaier, "A multi-perspective architecture for semantic code search," in *ACL*, 2020.
- [19] J. Gu, Z. Chen, and M. Monperrus, "Multimodal representation for neural code search," in *IEEE International Conference on Software Maintenance and Evolution, ICSME 2021, Luxembourg, September 27 - October 1, 2021*. IEEE, 2021, pp. 483–494. [Online]. Available: <https://doi.org/10.1109/ICSME52107.2021.00049>
- [20] Y. Wan, J. Shu, Y. Sui, G. Xu, Z. Zhao, J. Wu, and P. S. Yu, "Multi-modal attention network learning for semantic source code retrieval," in *ASE*. IEEE, 2019, pp. 13–25.
- [21] C. Ling, Z. Lin, Y. Zou, and B. Xie, "Adaptive deep code search," in *ICPC '20: 28th International Conference on Program Comprehension, Seoul, Republic of Korea, July 13-15, 2020*. ACM, 2020, pp. 48–59. [Online]. Available: <https://doi.org/10.1145/3387904.3389278>
- [22] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, S. Liu, L. Zhou, N. Duan, A. Svyatkovskiy, S. Fu, M. Tufano, S. K. Deng, C. B. Clement, D. Drain, N. Sundaresan, J. Yin, D. Jiang, and M. Zhou, "Graphcodebert: Pre-training code representations with data flow," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [Online]. Available: <https://openreview.net/forum?id=jLoC4ez43PZ>
- [23] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, "Codebert: A pre-trained model for programming and natural languages," in *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, ser. Findings of ACL, T. Cohn, Y. He, and Y. Liu, Eds., vol. EMNLP 2020. Association for Computational Linguistics, 2020, pp. 1536–1547. [Online]. Available: <https://doi.org/10.18653/v1/2020.findings-emnlp.139>
- [24] D. Guo, S. Lu, N. Duan, Y. Wang, M. Zhou, and J. Yin, "Unixcoder: Unified cross-modal pre-training for code representation," in *ACL (1)*. Association for Computational Linguistics, 2022, pp. 7212–7225.
- [25] Y. Wang, W. Wang, S. R. Joty, and S. C. H. Hoi, "Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation," in *EMNLP (1)*. Association for Computational Linguistics, 2021, pp. 8696–8708.
- [26] W. U. Ahmad, S. Chakraborty, B. Ray, and K. Chang, "Unified pre-training for program understanding and generation," in *NAACL-HLT*. Association for Computational Linguistics, 2021, pp. 2655–2668.
- [27] P. Jain, A. Jain, T. Zhang, P. Abbeel, J. Gonzalez, and I. Stoica, "Contrastive code representation learning," in *EMNLP (1)*. Association for Computational Linguistics, 2021, pp. 5954–5971.
- [28] N. D. Q. Bui, Y. Yu, and L. Jiang, "Self-supervised contrastive learning for code retrieval and summarization via semantic-preserving transformations," in *SIGIR*. ACM, 2021, pp. 511–521.
- [29] H. Husain, H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "Codesearchnet challenge: Evaluating the state of semantic code search," *CoRR*, vol. abs/1909.09436, 2019. [Online]. Available: <http://arxiv.org/abs/1909.09436>
- [30] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. Computer Vision Foundation / IEEE, 2020, pp. 9726–9735. [Online]. Available: <https://doi.org/10.1109/CVPR42600.2020.00975>
- [31] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," *CoRR*, vol. abs/1907.11692, 2019.
- [32] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer, "Summarizing source code using a neural attention model," in *ACL*, vol. 1. The Association for Computational Linguistics, 2016.
- [33] J. Zhang, X. Wang, H. Zhang, H. Sun, and X. Liu, "Retrieval-based neural source code summarization," in *ICSE*. IEEE / ACM, 2020.
- [34] A. LeClair, S. Jiang, and C. McMillan, "A neural model for generating natural language summaries of program subroutines," in *ICSE*, 2019, pp. 795–806.
- [35] E. Shi, Y. Wang, L. Du, H. Zhang, S. Han, D. Zhang, and H. Sun, "Cast: Enhancing code summarization with hierarchical splitting and reconstruction of abstract syntax trees," in *EMNLP*, 2021.
- [36] E. Shi, Y. Wang, L. Du, J. Chen, S. Han, H. Zhang, D. Zhang, and H. Sun, "On the evaluation of neural code summarization," in *ICSE*, 2022.
- [37] L. Du, X. Shi, Y. Wang, E. Shi, S. Han, and D. Zhang, "Is a single model enough? mucos: A multi-model ensemble learning approach for semantic code search," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 2994–2998.
- [38] W. Gu, Y. Wang, L. Du, H. Zhang, S. Han, D. Zhang, and M. R. Lyu, "Accelerating code search with deep hashing and code classification," in *ACL*, 2022.
- [39] A. Svyatkovskiy, S. Lee, A. Hadjitofi, M. Riechert, J. Franco, and M. Allamanis, "Fast and memory-efficient neural code completion," *Arxiv Preprint*, 2020. [Online]. Available: <https://arxiv.org/abs/1611.08307>
- [40] S. Proksch, J. Lerch, and M. Mezini, "Intelligent code completion with bayesian networks," *ACM Trans. Softw. Eng. Methodol.*, vol. 25, no. 1, pp. 3:1–3:31, 2015.
- [41] V. Raychev, M. T. Vechev, and E. Yahav, "Code completion with statistical language models," in *PLDI*, 2014, pp. 419–428.
- [42] M. Bruch, M. Monperrus, and M. Mezini, "Learning from examples to improve code completion systems," in *ESEC/FSE*, 2009, pp. 213–222.
- [43] Y. Wang and H. Li, "Code completion by modeling flattened abstract syntax trees as graphs," in *AAAI*, 2021.
- [44] W. Tao, Y. Wang, E. Shi, L. Du, S. Han, H. Zhang, D. Zhang, and W. Zhang, "A large-scale empirical study of commit message generation: models, datasets and evaluation," *Empirical Software Engineering*, vol. 27, no. 7, p. 198, 2022.
- [45] M. L. Vásquez, L. F. Cortes-Coy, J. Aponte, and D. Poshyvanyk, "Changescape: A tool for automatically generating commit messages," in *ICSE (2)*. IEEE Computer Society, 2015, pp. 709–712.
- [46] S. Xu, Y. Yao, F. Xu, T. Gu, H. Tong, and J. Lu, "Commit message generation for source code changes," in *IJCAI*. ijcai.org, 2019, pp. 3975–3981.
- [47] W. Tao, Y. Wang, E. Shi, L. Du, S. Han, H. Zhang, D. Zhang, and W. Zhang, "On the evaluation of commit message generation models: An experimental study," in *2021 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 2021, pp. 126–136.
- [48] E. Shi, Y. Wang, W. Tao, L. Du, H. Zhang, S. Han, D. Zhang, and H. Sun, "RACE: Retrieval-augmented commit message generation," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 5520–5530. [Online]. Available: <https://aclanthology.org/2022.emnlp-main.372>
- [49] C. Niu, C. Li, V. Ng, J. Ge, L. Huang, and B. Luo, "Spt-code: Sequence-to-sequence pre-training for learning source code representations," in *ICSE*. ACM, 2022, pp. 1–13.
- [50] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2006), 17-22 June 2006, New York, NY, USA*. IEEE Computer Society, 2006, pp. 1735–1742. [Online]. Available: <https://doi.org/10.1109/CVPR.2006.100>
- [51] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1597–1607.

- [52] T. Gao, X. Yao, and D. Chen, “Simcse: Simple contrastive learning of sentence embeddings,” in *EMNLP (1)*. Association for Computational Linguistics, 2021, pp. 6894–6910.
- [53] H. Fang and P. Xie, “CERT: contrastive self-supervised learning for language understanding,” *CoRR*, vol. abs/2005.12766, 2020.
- [54] J. M. Giorgi, O. Nitski, B. Wang, and G. D. Bader, “Declutr: Deep contrastive learning for unsupervised textual representations,” in *ACL/IJCNLP (1)*. Association for Computational Linguistics, 2021, pp. 879–895.
- [55] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” in *ICLR (Poster)*. OpenReview.net, 2018.
- [56] T. Devries and G. W. Taylor, “Improved regularization of convolutional neural networks with cutout,” *CoRR*, vol. abs/1708.04552, 2017.
- [57] A. G. Howard, “Some improvements on deep convolutional neural network based image classification,” in *ICLR (Poster)*, 2014.
- [58] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *ICML*, ser. JMLR Workshop and Conference Proceedings, vol. 37. JMLR.org, 2015, pp. 448–456.
- [59] X. Wang, Y. Wang, F. Mi, P. Zhou, Y. Wan, X. Liu, L. Li, H. Wu, J. Liu, and X. Jiang, “Syncobert: Syntax-guided multi-modal contrastive pre-training for code representation,” *arXiv preprint arXiv:2108.04556*, 2021.
- [60] Y. Ding, L. Buratti, S. Pujar, A. Morari, B. Ray, and S. Chakraborty, “Contrastive learning for source code with structural and functional properties,” *CoRR*, vol. abs/2110.03868, 2021.
- [61] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NIPS*, 2017, pp. 5998–6008.
- [62] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *CVPR*. Computer Vision Foundation / IEEE Computer Society, 2018, pp. 3733–3742.
- [63] Nvidia, “Nvidia a100 80gb pcie gpu,” *Product Brief*, 2022. [Online]. Available: [https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/PB-10577-001\\_v02.pdf](https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/PB-10577-001_v02.pdf)
- [64] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” 2018.
- [65] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *ICLR*, 2019.
- [66] Y. Wang, L. Du, E. Shi, Y. Hu, S. Han, and D. Zhang, “Cocogum: Contextual code summarization with multi-relational gnn on umls,” Microsoft, MSR-TR-2020-16. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/cocogum-contextual-code-summarization-with-multi-relational-gnn-on-umls>, Tech. Rep., 2020.
- [67] W. Gu, Z. Li, C. Gao, C. Wang, H. Zhang, Z. Xu, and M. R. Lyu, “Cradle: Deep code retrieval based on semantic dependency learning,” *Neural Networks*, vol. 141, pp. 385–394, 2021.
- [68] J. Huang, D. Tang, L. Shou, M. Gong, K. Xu, D. Jiang, M. Zhou, and N. Duan, “Cosqa: 20, 000+ web queries for code search and question answering,” in *ACL*, 2021.
- [69] H. Schütze, C. D. Manning, and P. Raghavan, *Introduction to information retrieval*. Cambridge University Press Cambridge, 2008, vol. 39.
- [70] S. E. Robertson and K. S. Jones, “Relevance weighting of search terms,” *Journal of the American Society for Information science*, pp. 129–146, 1976.
- [71] P. Jaccard, “Étude comparative de la distribution florale dans une portion des alpes et des jura,” *Bull Soc Vaudoise Sci Nat*, pp. 547–579, 1901.
- [72] M. Iyyer, V. Manjunatha, J. L. Boyd-Graber, and H. D. III, “Deep unordered composition rivals syntactic methods for text classification,” in *ACL (1)*. The Association for Computer Linguistics, 2015, pp. 1681–1691.
- [73] Y. Kim, “Convolutional neural networks for sentence classification,” in *EMNLP*. ACL, 2014, pp. 1746–1751.
- [74] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” in *EMNLP*. ACL, 2014, pp. 1724–1734.
- [75] C. Liu, X. Xia, D. Lo, C. Gao, X. Yang, and J. C. Grundy, “Opportunities and challenges in code search tools,” *ACM Comput. Surv.*, vol. 54, no. 9, pp. 196:1–196:40, 2022.
- [76] CoCoSoDa, “Replication package,” *ICSE*, 2023. [Online]. Available: <https://github.com/DeepSoftwareAnalytics/CoCoSoDa>
- [77] M. Mosbach, M. Andriushchenko, and D. Klakow, “On the stability of fine-tuning BERT: misconceptions, explanations, and strong baselines,” in *ICLR*. OpenReview.net, 2021.
- [78] T. Wang and P. Isola, “Understanding contrastive representation learning through alignment and uniformity on the hypersphere,” in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 9929–9939.
- [79] F. Wang and H. Liu, “Understanding the behaviour of contrastive loss,” in *CVPR*. Computer Vision Foundation / IEEE, 2021, pp. 2495–2504.
- [80] Y. Meng, C. Xiong, P. Bajaj, S. Tiwary, P. Bennett, J. Han, and X. Song, “COCO-LM: correcting and contrasting text sequences for language model pretraining,” in *NeurIPS*, 2021, pp. 23 102–23 114.
- [81] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [82] B. Li, H. Zhou, J. He, M. Wang, Y. Yang, and L. Li, “On the sentence embeddings from pre-trained language models,” in *EMNLP (1)*. Association for Computational Linguistics, 2020, pp. 9119–9130.
- [83] J. Gao, D. He, X. Tan, T. Qin, L. Wang, and T. Liu, “Representation degeneration problem in training natural language generation models,” in *ICLR (Poster)*. OpenReview.net, 2019.
- [84] A. Bansal, S. Haque, and C. McMillan, “Project-level encoding for neural source code summarization of subroutines,” in *ICPC*. IEEE, 2021, pp. 253–264.
- [85] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” in *NeurIPS*, 2020.