

Singapore Management University

## Institutional Knowledge at Singapore Management University

---

Research Collection School Of Computing and  
Information Systems

School of Computing and Information Systems

---

7-2021

### Self-supervised contrastive learning for code retrieval and summarization via semantic-preserving transformations

Duy Quoc Nghi BUI

Singapore Management University, dqnbui@smu.edu.sg

Yijun Yu

Lingxiao JIANG

Singapore Management University, lxjiang@smu.edu.sg

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)



Part of the [Software Engineering Commons](#)

---

#### Citation

BUI, Duy Quoc Nghi; Yijun Yu; and JIANG, Lingxiao. Self-supervised contrastive learning for code retrieval and summarization via semantic-preserving transformations. (2021). *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Conference, July 11–15*. 1-11.

Available at: [https://ink.library.smu.edu.sg/sis\\_research/6719](https://ink.library.smu.edu.sg/sis_research/6719)

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [cherylids@smu.edu.sg](mailto:cherylids@smu.edu.sg).

# Self-Supervised Contrastive Learning for Code Retrieval and Summarization via Semantic-Preserving Transformations

Nghi D. Q. Bui\*  
Trustworthy Software Engineering &  
Open Source Lab  
Huawei Ireland Research Center  
nghi.bui@huawei.com

Yijun Yu  
Trustworthy Software Engineering &  
Open Source Lab  
Huawei Ireland Research Center &  
The Open University, UK  
yijun.yu@huawei.com

Lingxiao Jiang  
School of Computing &  
Information Systems  
Singapore Management University  
lxjiang@smu.edu.sg

## Abstract

We propose *Corder*, a self-supervised contrastive learning framework for source code model. *Corder* is designed to alleviate the need of labeled data for code retrieval and code summarization tasks. The pre-trained model of *Corder* can be used in two ways: (1) it can produce vector representation of code which can be applied to code retrieval tasks that do not have labeled data; (2) it can be used in a fine-tuning process for tasks that might still require label data such as code summarization. The key innovation is that we train the source code model by asking it to recognize similar and dissimilar code snippets through a *contrastive learning objective*. To do so, we use a set of semantic-preserving transformation operators to generate code snippets that are syntactically diverse but semantically equivalent. Through extensive experiments, we have shown that the code models pretrained by *Corder* substantially outperform the other baselines for code-to-code retrieval, text-to-code retrieval, and code-to-text summarization tasks.

## CCS Concepts

• **Software and its engineering** → **Software libraries and repositories**; • **Information systems** → **Information retrieval**;

### ACM Reference Format:

Nghi D. Q. Bui, Yijun Yu, and Lingxiao Jiang. 2021. Self-Supervised Contrastive Learning for Code Retrieval and Summarization via Semantic-Preserving Transformations. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*, July 11–15, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3404835.3462840>

## 1 Introduction

Deep learning models for code have been found useful in many software engineering tasks, such as predicting bugs [30, 48, 72, 78], translating programs [16, 28], classifying program functionality [21, 56], searching code [27, 41, 60], generating comments from code

This work was mostly done when the first author was working in the School of Computing and Information Systems, Singapore Management University.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*SIGIR '21, July 11–15, 2021, Virtual Event, Canada*

© 2021 Association for Computing Machinery.  
ACM ISBN 978-1-4503-8037-9/21/07...\$15.00  
<https://doi.org/10.1145/3404835.3462840>

[6, 32, 68, 70], etc. These tasks can be seen as code retrieval where code could be either the documents to be found or the query to search. To build indices or models of source code for retrieval, a key step is to extract patterns of semantically equivalent or non-equivalent code snippets in large quantities. This is challenging for tremendous human efforts to collect and label code snippets. To overcome this challenge, one can depend on heuristics to label the code snippets automatically, e.g. by using test cases to compare programs [53]. The downside of such an approach is the extra costs associated with code execution, which may not always be possible. Another way is to collect free source code on hosting platforms such as Github, and extract the snippets that share similar comments [32], method names [6], or code documents, and then treat such snippets as semantically equivalent [33]. The drawback to this heuristic is that it can add a lot of noise because not all code snippets of identical comments, method names, or documents are indeed semantically equivalent. For example, Kang et al. [39] show that when pre-trained specifically for the method-name prediction task, the pre-trained Code2vec [6] model does not perform well for other code modeling tasks. Jiang et al. [38] perform further analysis to show the reason that methods with similar names are not necessarily semantically equivalent, which explains the poor transfer learning results of Kang et al. [39] on Code2vec since the model is forced to learn incorrect patterns of code.

To address these limitations, we develop *Corder*, a contrastive learning framework for code model that trains neural networks to identify semantically equivalent code snippets from a large set of transformed code. Essentially, the main goal of *Corder* is to invent a *pretext task* that enables the training of neural networks without the need for imprecise heuristics for identifying semantically equivalent programs. The pretext task that we implement here is called *instance discrimination*, which is related to recent work of Chen et al. [15]: a neural network is asked to discriminate instances of code snippets that are semantically equivalent from the instances of code snippets that are dissimilar. In the end, the neural model is trained with the knowledge of how different instances of code snippets should be close to each other in a vector space depending on how likely they are semantically equivalent. The key idea is to leverage program transformation techniques that transform a code snippet into different versions of itself. Although syntactically different from the originals, these transformed programs are semantically equivalent. Figure 1 shows an example of the transformed programs: Figure 1b shows a snippet semantically equivalent to the snippet in Figure 1a, with variables renamed. The snippet in

```

void insertionSort(int arr[]) {
  int n = arr.length;
  for (int i = 1; i < n; ++i) {
    int key = arr[i];
    int j = i - 1;
    while (j >= 0 && arr[j] > key){
      arr[j + 1] = arr[j];
      j = j - 1;
    }
    arr[j + 1] = key;
  }
}
a) Original Program

void insertionSort(int a[]) {
  int len = a.length;
  for (int k = 1; k < len; ++k) {
    int first = a[k];
    int j = k - 1;
    while (j >= 0 && a[j] > first){
      a[j + 1] = a[j];
      j = j - 1;
    }
    a[j + 1] = first;
  }
}
b) Same Program with Different Variable Names

void insertionSort(int arr[]) {
  int n = arr.length;
  for (int i = 1; i < n; ++i) {
    int key = arr[i];
    int j = i - 1;
    while (j >= 0 && arr[j] > key){
      arr[j + 1] = arr[j];
      j = j - 1;
    }
    arr[j + 1] = key;
  }
}
c) Same Program with Two Swapped Statements

```

**Figure 1: An Example of Semantically Equivalent Programs**

Figure 1c is another transformed version of Figure 1a, with two independent statements swapped. The goal is then to train the neural network with these semantically equivalent snippets and ensure they are embedded closely in the vector space.

Corder uses the *contrastive learning* methods that have been used in the self-supervised learning setting. The objective of contrastive learning is to simultaneously maximize the agreement between the differently transformed snippets of the same original snippet and minimize the agreement between the transformed snippets of different snippets. Updating the parameters of a neural network using this contrastive learning objective causes the representations of semantically equivalent snippets to be close to each other, while representations of non-similar snippets to be far apart. Once the model has been trained on our pretext task with the contrastive learning objective<sup>1</sup>, it can be used in two ways. First, since it has trained a neural network encoder (which is a part of the end-to-end learning process and can be instantiated with different encoders) and can be used to produce the representations of any source code. The vector representations of source code can be useful in many ways for code retrieval. Second, the pre-trained model can be fine-tuned with a small amount of labelled data to achieve good performance for other downstream tasks, such as code summarization. In this work, we consider three concrete tasks that can leverage our pre-trained model, namely: code-to-code retrieval, text-to-code retrieval, and code summarization (which is code-to-text). We have trained different Corder instances with different encoders on large-scale Java datasets to evaluate their effectiveness for the tasks.

To summarize, our major contributions are as follows:

- We explore a novel perspective of learning source code models from unlabeled data. Unlike existing work that uses imprecise heuristics to produce labelled data for training, we adapt program transformation techniques to generate precise semantically-equivalent code snippets for training. To the best of our knowledge, we are the first to use the program transformation technique for self-supervised learning of source code models.
- We develop Corder, a self-supervised contrastive learning framework, to identify semantically-equivalent code snippets generated from program transformation operators.
- We conduct extensive evaluations to demonstrate that the Corder pretext task is better than others in learning source code models in two ways: (1) we use the pre-trained models to produce vector representations of code and apply such representations in the *unsupervised* code-to-code retrieval task. The results show that any neural network encoder trained on the Corder pretext task outperforms the same encoders trained on other pretext

<sup>1</sup>We call this the Corder pretext task.

tasks with a significant margin. Moreover, our technique outperforms other baselines that were designed specifically for code-to-code retrieval, such as FaCoy [41] significantly; (2) we use the pre-trained models in a fine-tuning process for *supervised* code modeling tasks, such as text-to-code retrieval and code summarization. The results show that our pre-trained models on the Corder pretext task perform better than training the code models from scratch and other pretext tasks, by a large margin.

## 2 Related Work

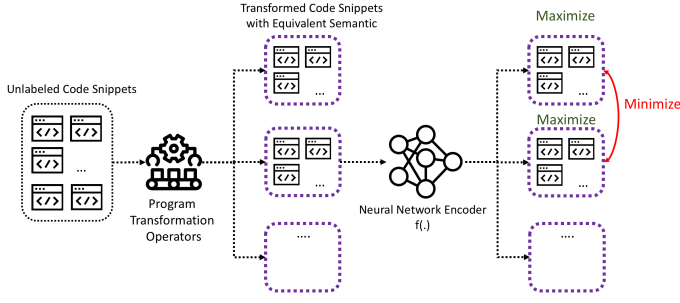
*Self-Supervised Learning* has made tremendous strides in the field of visual learning [25, 26, 40, 46, 52, 76], and for quite some time in the field of natural language processing [22, 42, 45, 47, 54, 62, 69, 71]. Such techniques allow for neural network training without the need for human labels. Typically, a self-supervised learning technique reformulates an unsupervised learning problem as one that is supervised by *generating virtual labels automatically from existing (unlabeled) data*. *Contrastive learning* has emerged as a new paradigm that brings together multiple forms of supervised learning problem as the task to compare similar and dissimilar items, such as siamese neural networks [9], contrastive predictive coding [57], triplet loss [61]. Contrastive learning methods specifically minimize the distance between similar data (positives) representations and maximize the distance between dissimilar data (negatives).

*Deep Learning Models of Code* : There has been a huge interest in applying deep learning techniques for software engineering tasks such as program functionality classification [10, 11, 55, 75], bug localization [18, 29, 37, 58], code summarization [2, 24, 66], code clone detection [13, 75], program refactoring [32], program translation [12, 16], and code synthesis [5, 8]. Allamanis et al. [3] extend ASTs to graphs by adding a variety of code dependencies as edges among tree nodes, intended to represent code semantics, and apply Gated Graph Neural Networks (GGNN) [49] to learn the graphs; Code2vec [6], Code2seq [4], and ASTNN [75] are designed based on splitting ASTs into smaller ones, either as a bag of path-contexts or as flattened subtrees representing individual statements. They use various kinds of Recurrent Neural Network (RNN) to learn such code representations. Surveys on code embeddings [17, 35] present evidence to show that there is a **strong need to alleviate the requirement of labeled data for code modeling** and encourage the community to invest more effort into the methods of learning source code with unlabeled data. Unfortunately, there is little effort towards designing the source code model with unlabeled data: Yasunaga and Liang [73] presents a self-supervised learning paradigm for program repair, but it is designed specifically for program repair only. There are methods, such as [23, 34] that perform pretraining source code data on natural language model (BERT, RNN, LSTM), but they simply train the code tokens similar to the way pretrained language models on text do, so they miss a lot of information about syntactical and semantic features of code that could have been extracted from program analysis.

## 3 Approach

### 3.1 Approach Overview

Figure 2 presents an overview of our approach. Overall, this framework comprises the following three major components.



**Figure 2: Overview of Corder pretext task. Unlabeled code snippets from a large codebase go through a program transformation module. Snippets in the purple dashed box are transformed snippets from the same original snippet. The goal is to maximize the similarity of the snippets in the same purple dashed box and minimize the similarity of snippets across different boxes**

---

**Algorithm 1** Corder’s learning algorithm

---

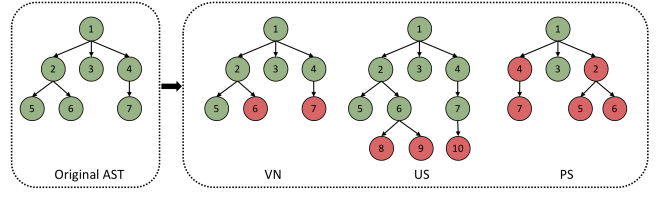
```

1: input: batch size  $N$ , encoder  $f$ , set of transformation operators  $\mathcal{T}$ .
2: for sampled minibatch  $\{p_k\}_{k=1}^N$  do
3:   for all  $k \in \{1, \dots, N\}$  do
4:     draw two transformation operators  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
5:     # the first transformation
6:      $\tilde{p}_i = t(p_k)$ 
7:      $v_i = f(\tilde{p}_i)$ 
8:     # the second transformation
9:      $\tilde{p}_j = t'(p_k)$ 
10:     $v_j = f(\tilde{p}_j)$ 
11:  end for
12:  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
13:     $s_{i,j} = z_i^\top z_j / (\|z_i\| \|z_j\|)$  # pairwise similarity
14:  end for
15:  define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j})}{\sum_{k=1}^{2N} \mathbf{1}_{k \neq i} \exp(s_{i,k})}$ 
16:   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
17:  update networks  $f$  to minimize  $\mathcal{L}$ 
18: end for
19: return encoder network  $f(\cdot)$ 

```

---

- A **program transformation module** that transforms a given code snippet  $p$ , resulting in two transformed programs of the code snippets, denoted  $\tilde{p}_i$  and  $\tilde{p}_j$ .
- A **neural network encoder**  $f(\cdot)$  that receives an intermediate representation of a code snippet (such as Abstract Syntax Tree (AST)) and map it into a vector representation. In our case, it should map  $\tilde{p}_i$  and  $\tilde{p}_j$  into two code vectors  $v_i$  and  $v_j$ , respectively.
- A **contrastive loss function** is defined for the contrastive learning task. Given a set  $\{\tilde{p}_k\}$  containing a positive (semantically similar) pair of examples  $p_i$  and  $\tilde{p}_j$ , the *contrastive prediction task* aims to identify  $\tilde{p}_j$  in  $\{\tilde{p}_k\}_{k \neq i}$  for a given  $p_i$ .



**Figure 3: Example of how the AST structure is changed with different transformation operators**

### 3.2 Approach Details

With the above components, here we describe the Corder training process in the following steps. Also, a summarization of the proposed algorithm is depicted in Algorithm 1.

- A mini-batch of  $N$  samples is randomly selected from a large set of code snippets. Each code snippet  $p$  in  $N$  is applied with two different randomly selected transformation operators, resulting in  $2N$  transformed code snippets:  $\tilde{p}_i = t(p)$ ;  $\tilde{p}_j = t'(p)$ ;  $t, t' \sim \mathcal{T}$ , where  $p$  is the original code snippet,  $\tilde{p}_i$  and  $\tilde{p}_j$  are transformed code snippets by applying two transformation operators  $t$  and  $t'$  into  $p$ , respectively.  $t$  and  $t'$  are randomly chosen from a set of available operators  $\mathcal{T}$ .
- Each of the transformed snippet  $\tilde{p}_i$  and  $\tilde{p}_j$  will be fed into the same encoder  $f(\cdot)$  to get the embedding representations:  $v_i = f(\tilde{p}_i)$ ;  $v_j = f(\tilde{p}_j)$ .
- We use the Noise Contrastive Estimate (NCE) loss function [15] to compute the loss. Let  $\text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$  denote the dot product between  $\ell_2$  normalized  $\mathbf{u}$  and  $\mathbf{v}$  (i.e. cosine similarity). Then the loss function for a pair of representations  $(v_i, v_j)$  is defined as  $\ell(i, j) = -\log \frac{\exp(\text{sim}(v_i, v_j))}{\sum_{k=1}^{2N} \mathbf{1}_{k \neq i} \exp(\text{sim}(v_i, v_k))}$ , where  $\mathbf{1}_{k \neq a} \in \{0, 1\}$  is an indicator function evaluating to 1 iff  $k \neq i$ . Noted that for a given positive pair, the other  $2(N-1)$  transformed code snippets are treated as negative samples. We calculate the loss for the same pair a second time as well where the positions of the samples are interchanged. The final loss is computed across all pairs in a mini-batch can be written as:  $L = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ .

**3.2.1 Program Transformation Operators.** Our key idea to enable the encoder to learn a set of diverse code features without the need for labeled data is to generate multiple versions of a program without changing its semantics. To do so, we apply a set of semantic-preserving program transformation operators to generate such different variants. There are many methods for transforming the code [59], and the more sophisticated a transformation is, in principle, the better the encoder can learn the essential semantic of the code. In this work, we mainly apply the following transformations: variable renaming, adding dead code (unused statements), permuting statements, loop exchange, and switch to if, to reflect different ways to change the structure of the Abstract Syntax Tree (AST). We evaluate how different transformations can effect the performance of the encoder in Section 6.1.

- **Variable Renaming (VN)** is a refactoring method that renames a variable, where the new name of the variable is taken randomly from a set of variable vocabulary in the training set. Noted that each time this operator is applied to the same program, the

variable names are renamed differently. This operator does not change the structure of the AST representation of the code, it only changes the textual information, which is a feature of a node in the AST.

- **Unused Statement (US)** is to insert dead code fragments, such as unused statement(s) to a randomly selected basic block in the code. We traverse the AST to identify the blocks and randomly select one block to insert predefined dead code fragments into it. This operator will add more nodes to the AST. To diversify the set of transformed programs, we prepare a large set of unused statement(s). When the operator is applied, random statements in the set is selected to added into the code block, i.e., a transformed snippet is different each time we apply the same operator.
- **Permutation of Statements (PS)** is to swap two statements that have no dependency on each other in a basic block in the code. We traverse the AST and analyze the data dependency to extract all of the possible pairs of swap-able statements. If a program only contains one such pair, it will generate the same output every time we apply the operator, otherwise, the output will be different.
- **Loop Exchange (LX)** replaces for loops with while loops or vice versa. We traverse the AST to identify the node the defines the for loop (or the while loop) then replace one with another with modifications on the initialization, the condition, and the afterthought.
- **Switch to If (SF)** replaces a switch statement in the method with its equivalent if statement. We traverse the AST to identify a switch statement, then extract the subtree of each case statement of the switch and assign it to a new if statement.

Each of the transformation operators above is designed to change the structure representation of the source code differently. For example, with VR, we want the NN to understand that even the change in textual information does not affect the semantic meaning of the source code, inspired by a recent finding of Zhang et al. [74]. It is suggested that the source code model should be equipped with adversarial examples of token changes to make the model become more robust. With US, we want the NN still to learn how to catch the similarity between two similar programs even though the number of nodes in the tree structure has increased. With PS, the operator does not add nodes into the AST but it will alter the location of the subtrees in the AST, we want the NN to be able to detect the two similar trees even if the positions of the subtrees have changed. Figure 3 illustrates how the AST structure changes with the corresponding transformation operator.

**3.2.2 Neural Network Encoder for Source Code** The neural network can also be called as an *encoder*, written as a function  $f(\cdot)$ . The encoder receives the intermediate representation (IR) of code and maps it into a code vector embedding  $\vec{v}$  (usually a combination of various kinds of code elements), then  $\vec{v}$  can be fed into the next layer(s) of a learning system and trained for an objective function of the specific task of the learning system. The choice of the encoder depends mostly on the task and we will rely on previous work to choose suitable encoders for a particular task, which will be presented in Section 5.

## 4 Use Cases

We present three tasks (code-to-code retrieval, text-to-code retrieval, and code summarization) to make good uses of the pre-trained Corder models in two ways.

### 4.1 Using the Pre-trained Encoders to Produce Code Vectors for Downstream Task

The first way to use pre-trained encoders from our Corder pretext task is to use such encoders to produce the vector representations of code. Then the representations can be applicable for a downstream task, such as code-to-code retrieval.

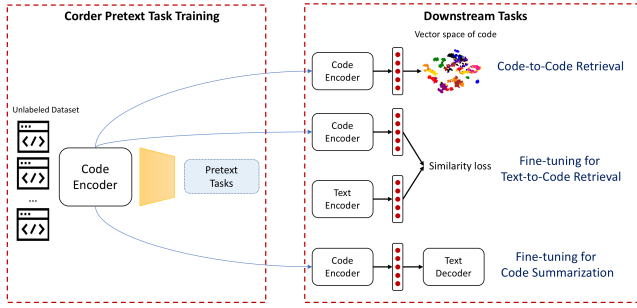
**4.1.1 Code-to-Code Retrieval** *Code-to-code* search is useful for developers to find other code in a large codebase that is similar to a given code query. Most of the work that is designed for code-to-code retrieval, such as Facoy [41], Krugle [1] is based on the simple text mining approach or traditional code clone detection method. These techniques required tremendous effort of handcraft feature engineering to extract good features of code. In our case, we adapt pre-trained source code encoders from the Corder pretext task to map any code snippet into a vector representation, then we perform the retrieval task based on the vectors (see Figure 4, Code-to-Code Retrieval). Assume that we have a large codebase of snippets, we used the pre-trained encoders to map the whole codebase into representations. Then for a given code snippet as a query, we map such query into vector representation too. Then, one can find the top-k nearest neighbors of such query in the vector space, using cosine similarity as the distance metric, and finally can retrieve the list of candidate snippets. These snippets are supposed to be semantical equivalent to the query.

### 4.2 Fine-Tuning the Encoders for Supervised Learning Downstream Tasks

A paradigm to make good use of a large amount of unlabeled data is self-supervised pre-training followed by a supervised fine-tuning [15, 31], which reuses parts (or all) of a trained neural network on a certain task and continue to train it or simply using the embedding output for other tasks. Such fine-tuning processes usually have the benefits of (1) speeding up the training as one does not need to train the model from randomly initialized weights and (2) improving the generalizability of the downstream model even when there are only small datasets with labels. As shown in Figure 4, the encoder is used as a pre-trained model in which the weights resulting from the Corder pretext task are transferred to initialize the model of the downstream supervised learning tasks.

**4.2.1 Text-to-Code Retrieval** This task is to, given a natural language as the query, the objective is to find the most semantically related code snippets from a collection of codes [27, 33]. Note that this is different from the *code-to-code* retrieval problem, in which the query is a code snippet. The deep learning framework used in the literature for this task is to construct a bilateral neural network structure, which consists of two encoders, one is a *natural language encoder* (such as BERT, RNN, LSTM) to encode text into text embedding, the other is a *source code encoder* to encode an immediate source code representation into the code embedding [27, 33]. Then,





**Figure 4: Process on how Corder pre-trained model can be applied in different downstream tasks**

from text embedding and code embedding, a mapping function is used to push the text and the code to be similar to the vector space, called a shared embedding between the code and the text. In the retrieval process, the text description is given, and we use its embedding to retrieve all the embeddings of the code snippets that are closest to the text embedding. In the fine-tuning process, the source code encoder that has been pre-trained on the Corder pretext task will be used to initialize the source code encoder, the parameters of the text encoder will be initialized randomly.

**4.2.2 Code Summarization** The purpose of this task is to predict a concise text description of the functionality of the method given its source code [7]. Such descriptions typically appear as documentation of methods (e.g. "docstrings" in Python or "JavaDocs" in Java). This task can be modeled as a translation task where the aim is to translate a source code snippet into a sequence of text. As such, the encoder-decoder model, such as seq2seq [64] is usually used in the literature for this task. In our case, the encoder can be any code modeling technique, such as TBCNN [55], Code2vec [6], LSTM or Transformer. In the fine-tuning process, the source code encoder that has been pre-trained on the Corder pretext task will be used to initialize for the source code encoder.

## 5 Empirical Evaluation

### 5.1 Settings

**5.1.1 Data Preparation** As presented, we will perform the evaluation on three tasks, namely, code-to-code search, text-to-code search, and code summarization. We used the JavaSmall and JavaMed datasets that have been widely used recently for code modeling tasks [4, 6]. JavaSmall is a dataset of 11 relatively large Java projects from GitHub, which contains about 700k examples. JavaMed is a dataset of 1000 top-starred Java projects from GitHub which contains about 4M examples.

Then, we parse all the snippets into ASTs using SrcML [19]. We also perform the transformation on all of the ASTs to get the transformed ASTs based on the transformation operators described in Section 3.2.1, having the ASTs as well as the transformed ASTs. It should be noted that SrcML is a universal AST system, which means that it uses the same AST representations for multiple languages (Java, C#, C++, C). This enables the model training on each of the languages once and they can be used in other languages. Another

note is that the two datasets are not the ones used for evaluation purposes; they are only for the purpose of training the Corder pretext task on different encoders. We describe the evaluation datasets used for each of the tasks separately in each of the subsections.

**5.1.2 Encoders** We choose a few well-known AST-based code modeling techniques as the encoder  $f(\cdot)$ , which are Code2vec [6], TBCNN [55], We also include two token-based techniques by treating source code simply as sequences of tokens and using a neural machine translation (NMT) baseline, i.e. a 2-layer Bi-LSTM, and the Transformer [67]. A common setting used among all these techniques is that they all utilize both node type and token information to initialize a node in ASTs.

We set both the dimensionality of type embeddings and text embeddings to 128. Note that we try our best to make the baselines as strong as possible by choosing the hyper-parameters above as the "optimal settings" according to their papers or code. Specifically, for Code2vec [6]<sup>2</sup> and Code2seq [4]<sup>3</sup>, since Code2seq is a follow-up work of Code2vec (different in the decoder layer to predict the sequence), we follow the settings in Code2seq to set the size of each LSTM encoders for ASTs to 128 and the size of LSTM decoder to 320. We set the number of paths sampled for each AST to 200 as suggested, since increasing this parameter does not improve the performance. TBCNN [55] uses a tree-based convolutional layer with three weight matrices serving as model parameters to accumulate children's information to the parent, each will have the shape of  $128 \times 128$ . We set the number of convolutional steps to 8. For Transformer [67], we choose to set the number of layers to 5 and the attention dimension size to 128. Finally, for the 2-layer Bi-LSTM, we followed the strategy from Alon et al. [4] by assigning the token embedding size to 128, the size of the hidden unit in the encoder to 128, and the default hyperparameters of OpenNMT [43].

**5.1.3 Research Questions** We want to answer two research questions specifically through the evaluations:

- (1) Are the code vectors generated by the pre-trained models (with various encoders) useful in the unsupervised space searching task (code-to-code search in particular)?
- (2) Can the pre-trained models be used in a fine-tuning process to improve the performance of the downstream supervised models for text-to-code search and code summarization without training the models from scratch?

## 5.2 Using Pre-trained Encoders to Produce Code Representations for Code-to-Code Retrieval

**5.2.1 Datasets, Metrics, and Baselines** Given a code snippet as the input, the task aims to find the most semantically related code from a collection of candidate codes. The datasets we used to evaluate for this task are:

- OJ dataset [55] contains 52000 C programs with 104 classes, which results in 500 programs per class. Since the dataset is for C++, we translate the whole dataset with the C++ to Java Converter<sup>4</sup> to make the language of the evaluation dataset aligned with the

<sup>2</sup><https://github.com/tech-srl/code2vec>

<sup>3</sup><https://github.com/tech-srl/code2seq>

<sup>4</sup>[https://www.tangiblesoftware.com/product\\_details/cplusplus\\_to\\_java\\_converter\\_details.html](https://www.tangiblesoftware.com/product_details/cplusplus_to_java_converter_details.html)

pretrained models for Java (see Section 5.1). Then we use the data that has been translated to Java for evaluation.

- BigCloneBench (BCB) dataset [65] contains 25,000 Java projects, cover 10 functionalities and including 6,000,000 true clone pairs and 260,000 false clone pairs. This dataset has been widely used for code clone detection task.

The OJ and BigCloneBench datasets have been widely used for the code clone detection task. The code clone detection task is to detect semantically duplicated (or similar) code snippets in a large codebase. Thus these datasets are also suitable for the code-to-code retrieval task, with the aim to find the most semantically related codes given the code snippet as the query.

We randomly select 50 programs per class as the query, so that the total number of queries is 5200 for 104 classes. For each of the queries, we want to retrieve all of the semantically similar code snippets, which are the programs in the same class of the query. With OJ, each query can have multiple relevant results, so that we use Mean Average Precision (MAP) as the metric to evaluate for the code-to-code search on the OJ dataset. Mean average precision for a set of queries is the mean of the average precision scores for each query, which can be calculated as  $MAP = \frac{\sum_{q=1}^Q AveP(q)}{Q}$ , where  $Q$  is the number of queries in the set and  $AveP(q)$  is the average precision for a given query  $q$ .

For the BCB dataset, since the size of the dataset is large, we reduce the size by randomly select 50,000 sample clone pairs and 50,000 samples none clone pairs and evaluate within these pairs. Then within the clone pairs, we again randomly select 5000 pairs and pick one code snippet of a pair as the query. Let's denote a clone pair as  $p = (c_1, c_2)$ , we pick  $c_1$  as the query. For each of the query  $c_1$ , we want to retrieve the  $c_2$ , which is the snippet that is semantically identical to the query. With BCB, the assumption is that each query has only one relevant result so that we use Mean Reciprocal Rank (MRR) as the metric to evaluate for the task. Mean Reciprocal Rank is the average of the reciprocal ranks of results of a set of queries  $Q$ . The reciprocal rank of a query is the inverse of the rank of the first hit result. The higher the MRR value, the better the code search performance. MRR can be calculated as follows:  $MRR = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \frac{1}{rank_i}$ . Note that for both datasets, we limited the number of return results to 10. We use these baselines for the code-to-code retrieval task:

- Word2vec: the representation of the code snippet can be computed by simply calculate the average of the representations of all of the token in the snippet
- Doc2vec: we use Gensim<sup>5</sup> to train the Doc2vec model on the JavaMed dataset and use the method provided by Gensim to infer the representation for a code snippet
- ElasticSearch: we treat the code snippet as a sequence and use the text tokenizer provided by ElasticSearch to index the code token and use ElasticSearch as a fuzzy text search baseline.
- Facoy [41] is a search engine that is designed specifically for code-to-code search.

Besides the baselines above, we also want to see if our Corder pretext task performs better than the other pretext task for the same encoder. Among the encoders, the Transformer [67] can be

**Table 1: Results of code-to-code retrieval. For BigCloneBench (BCB), the metric is MAP. For OJ, the metric is MRR**

Model	Pre-training	Dataset	Performance	
			BCB (MRR)	OJ (MAP)
ElasticSearch	-	-	0.131	0.235
Word2Vec	-	JavaMed	0.255	0.212
Doc2Vec	-	JavaMed	0.289	0.334
FaCoy	-	JavaMed	0.514	0.585
Code2Vec	MNP	JavaSmall	0.374	0.529
	MNP	JavaMed	0.453	0.621
	Corder	JavaSmall	0.561	0.631
	Corder	JavaMed	<b>0.633</b>	<b>0.735</b>
TBCNN	Corder	JavaSmall	0.654	0.710
	Corder	JavaMed	<b>0.832</b>	<b>0.856</b>
Bi-LSTM	Corder	JavaSmall	0.620	0.534
	Corder	JavaMed	<b>0.742</b>	<b>0.690</b>
Transformer	Masked LM	JavaSmall	0.580	0.534
	Masked LM	JavaMed	0.719	0.693
	Corder	JavaSmall	0.640	0.720
	Corder	JavaMed	<b>0.825</b>	<b>0.841</b>

pre-trained with other pretext tasks, such as the masked language modeling, where a model uses the context words surrounding a [MASK] token to try to predict what the [MASK] word should be. Code2vec [6] is also applicable for another pretext task, which is the method name prediction (MNP). The path encoder in Code2vec can encode the method body of a code snippet, then use the representation of the method body to predict the method name. With this, the Code2vec model can be pre-trained with MNP as a pretext task. The path encoder of the Code2vec for the method name prediction task can be reused to produce representation for any code snippet. For such reasons, we include 2 additional baselines, which are a pre-trained Transformer on the masked language model on the JavaMed dataset, and a pre-trained Code2vec on MNP on the JavaMed dataset.

**5.2.2 Results** Table 1 shows the results of the code-to-code retrieval task. The column "Pre-training" with different options, such as "Corder", "Masked LM", "MNP" means that an encoder is trained with a different pretext task. ElasticSearch, Word2vec, Doc2vec, and FaCoy are not applicable for such pretext tasks, hence "-" is used for these techniques in the column "Pre-training". The column "Dataset" means that an encoder is trained on a specific dataset, such as JavaSmall and JavaMed.

As one can see, ElasticSearch, an information retrieval approach, performs worst among the baselines. Word2vec and Doc2vec perform better but the results are still not so good. Code2vec and Bi-LSTM, when pre-training with the Corder process on the JavaMed, can perform better than FaCoy, a method designed specifically for code-to-code retrieval. Code2vec, when pre-training with the method name prediction (MNP) pretext task, performs much worse than the pre-training with the Corder pretext task. Transformer, when pre-training with the masked language model (Masked-LM) pretext task, performs much worse than the pre-training with the Corder pretext task. This shows that our proposed pretext task performs better than the other pretext tasks to train the representation of the source code.

<sup>5</sup><https://github.com/RaRe-Technologies/gensim>

Table 2: Results of text-to-code retrieval

Model	Pre-training	Dataset	P@1	P@5	P@10	MRR
NBow	-	-	0.394	0.581	0.603	0.384
	-	-	0.406	0.529	0.564	0.395
Code2vec	MNP	JavaSmall	0.415	0.538	0.572	0.409
	MNP	JavaMed	0.435	0.546	0.583	0.420
	Corder	JavaSmall	0.512	0.578	0.610	0.446
	Corder	JavaMed	<b>0.549</b>	<b>0.608</b>	<b>0.625</b>	<b>0.592</b>
TBCNN	-	-	0.506	0.581	0.632	0.551
	Corder	JavaSmall	0.541	0.620	0.658	0.658
	Corder	JavaMed	<b>0.640</b>	<b>0.710</b>	<b>0.758</b>	<b>0.702</b>
Bi-LSTM	-	-	0.469	0.540	0.702	0.630
	Corder	JavaSmall	0.532	0.581	0.723	0.619
	Corder	JavaMed	<b>0.567</b>	<b>0.639</b>	<b>0.768</b>	<b>0.661</b>
Transformer	-	-	0.534	0.653	0.793	0.651
	Masked LM	JavaSmall	0.567	0.627	0.683	0.630
	Masked LM	JavaMed	0.632	0.710	0.753	0.672
	Corder	JavaSmall	0.604	0.698	0.845	0.687
	Corder	JavaMed	<b>0.662</b>	<b>0.756</b>	<b>0.881</b>	<b>0.728</b>

### 5.3 Fine-tuning Pre-trained Encoders for Text-to-code Retrieval

5.3.1 *Datasets, Metrics, and Baselines* Given a natural language as input, the task aims to find the most semantically related code from a collection of candidate codes. We use the dataset released by DeepCS [27], which consists of approximately 16 million pre-processed Java methods and their corresponding docstrings.

For the metrics, we use Precision at k (Precision@k) and Mean Reciprocal Rank to evaluate this task. Precision@k measures the percentage of relevant results in the top k returned results for each query. In our evaluations, it is calculated as follows:  $Precision@k = \frac{\#relevant\ results\ in\ the\ top\ k\ results}{k}$ . Precision@k is important because developers often inspect multiple results of different usages to learn from. A better code retrieval engine should allow developers to inspect less noisy results. The higher the metric values, the better the code search performance. We evaluate and Precision@k when the value of k is 1, 5, and 10. These values reflect the typical sizes of results that users would inspect.

We choose to use the three methods presented in CodeSearchNet [33] for the text-to-code retrieval models, which are: neural bag-of-words, 2-layer BiLSTM, and Transformer. We also include Tree-based CNN (TBCNN) [55] and Code2vec [6] which are AST-based encoders that receive the AST representation of the code snippets as the input. We perform evaluations under two settings: (1) train from scratch and (2) fine-tune with a pre-trained model. In the second setting, each of the encoders will be pre-trained through the Corder pretext task, then the pre-trained encoder will be used for the fine-tuning process. We also include the pre-trained Code2vec model from the method name prediction (MNP) task to demonstrate that our Corder pretext task is better in a fine-tuning process than the pre-trained model from the MNP task.

5.3.2 *Results* Table 2 shows the performance of text-to-code retrieval task. The column "Pre-training" with different options, such as "-", "MNP", "Corder", means that an encoder is pre-trained on different pretext tasks. "-" means that there is no pretext task applied for the model and the model is trained from scratch. The column "Dataset", with different options, such as "-", "JavaSmall", "JavaMed",

Table 3: Results of code summarization

Model	Pre-training	Dataset	BLEU
MOSES	-	-	11.56
IR	-	-	14.32
	-	-	22.89
Code2seq	MNP	JavaSmall	23.14
	MNP	JavaMed	24.20
	Corder	JavaSmall	24.23
	Corder	JavaMed	<b>26.56</b>
TBCNN	-	-	21.15
	Corder	JavaSmall	23.56
	Corder	JavaMed	<b>25.39</b>
Bi-LSTM	-	-	23.98
	Corder	JavaSmall	24.21
	Corder	JavaMed	<b>25.50</b>
Transformer	-	-	22.85
	Masked LM	JavaSmall	23.10
	Masked LM	JavaMed	24.78
	Corder	JavaSmall	25.69
	Corder	JavaMed	<b>26.41</b>

means that the encoder is pre-trained on a pretext task with a specific dataset. There are 3 observations : (1) Corder pre-training task on any of the model improves the performance significantly; (2) pre-training on a larger dataset improves the results with a higher margin than pre-training on a smaller dataset, and (3) Corder pretext task for Code2vec performs better than the MNP task to fine-tune the model for text-to-code retrieval.

### 5.4 Fine-tuning Pre-trained Encoders for Code Summarization

5.4.1 *Dataset, Metric, and Baselines* For this task, we consider predicting a full natural language sentence given a short code snippet. We also use the Java dataset provided by DeepCS [27], which consists of approximately 16 million preprocessed Java methods and their corresponding docstrings. The target sequence length in this task is about 12.3 on average. Since this dataset consists of a parallel corpus of code snippets and docstrings, it is suitable for either the text-to-code retrieval task or the code summarization task.

To measure the prediction performance, we follow [6] to use the BLEU score as the metric. For the baselines, we present results compared to 2-layer bidirectional LSTMs, Transformer, and Code2seq [4], a state-of-the-art model for code summarization task. We provide a fair comparison by splitting tokens into subtokens and replacing UNK during inference. We also include numbers from the baselines used by Iyer et al. [36], such as MOSES [44] and an IR-based approach that use Levenshtein distance to retrieve the description.

5.4.2 *Results* Table 3 shows the performance of Corder pretraining on the code summarization task. As seen, pre-training the model with the Corder pretext task outperform the other pre-training task, such as MNP or Masked LM in term of BLEU score with a significant margin for any of the encoder.

### 5.5 Compared Against Supervised Methods

An interesting question that one might ask is how Corder’s performance in comparison with some supervised methods? The results for the tasks set out in the previous subsections may not be appropriate to answer this question because the three tasks are only used to measure how well the Corder Pretext task performs in a retrieval



**Table 4: Comparison against supervised representation learning methods on code classification task.**

Methods	Accuracy
<i>Pretrained from Supervised Methods</i>	
AdaSent	0.64
InferSent	0.56
Code2vec	0.75
<i>Pretrained from Corder Pretext</i>	
Code2vec	0.82
TBCNN	0.85
Bi-LSTM	0.76
Transformer	0.80

task (code-to-code retrieval) or in a fine-tuning process (text-to-code retrieval and code summarization), it is not clear if the Corder Pretext task can perform better than other supervised learning methods specifically trained for the tasks. As such, a more appropriate *supervised* task is needed to answer this question. A well-known evaluation protocol widely used to measure if the self-supervised learning methods can beat the supervised ones in natural language processing is to train the classifier (usually the Logistic Regression) on top of the embeddings provided by *trained encoders on a supervised task* for classification tasks [42, 50], and Accuracy is used as a proxy for representation quality. We adopt the *code classification* task on the OJ dataset [55] and follow the similar evaluation protocol from Logeswaran and Lee [50] to evaluate if the embeddings provided by the self-supervised learning techniques are better than the supervised ones. The code classification task is to, given a piece of code, classify the functionality class it belongs to. The term *code classification* is sometimes used interchangeably with the *code summarization* task [63], in which we want to automatically assign a label (or description) to a code snippet; it is just the way in which the label is generated can be different.

We produce the embeddings from all of the encoders for all of the training samples, then we train a classifier on top of these embeddings. We adapt a few strong state-of-the-art learning techniques for sentence representations in NLP to model the source code as a sequence of tokens, which are AdaSent [77], InferSent [20]. The encoders from these two techniques are trained from supervised tasks. We also use the encoder from the pretrained Code2vec on the method name prediction task as another baseline. The embeddings produced by these techniques will also be used to train the multi-class classifier. We choose Logistic Regression as the classifier in this study. We use Accuracy as the metric to measure the performance of the code classification task. Table 4 shows the results of this analysis. The performance of the encoders trained from the Corder pretext performs significantly better than the other supervised learning methods.

## 5.6 Summary & Threats to Validity

Corder outperforms most of the baselines across three tasks: code-to-code retrieval, text-to-code retrieval, and code summarization. We also show that the embeddings produced by Corder perform better than the embeddings produced by the other supervised learning methods with code classification.

**Table 5: Results on Analysis on the Impact of Different Transformation Operators. TTC = Text-to-Code Retrieval (MRR as the metric), CS = Code Summarization (BLEU as the metric)**

Models	Ops	Original		Downsampled	
		TTC	CS	TTC	CS
Code2vec	VR	0.434	19.56	0.202	16.78
	US	0.498	21.45	0.345	19.06
	PS	0.552	24.22	0.385	20.01
	SF	0.401	19.11	0.401	19.11
	LX	0.423	21.43	0.320	20.18
	All	0.592	26.56	0.419	21.25
TBCNN	VR	0.421	20.11	0.246	17.89
	US	0.562	23.56	0.368	18.24
	PS	0.603	22.98	0.320	19.50
	LX	0.519	22.20	0.311	19.82
	SF	0.354	18.75	0.461	18.75
	All	0.702	25.39	0.398	21.05
Bi-LSTM	VR	0.423	21.53	0.302	17.33
	US	0.601	22.42	0.401	18.32
	PS	0.621	23.57	0.398	19.56
	LX	0.529	20.89	0.328	20.22
	SF	0.412	19.34	0.412	19.34
	All	0.661	25.50	0.435	20.84
Transformer	VR	0.411	20.57	0.286	17.82
	US	0.581	23.56	0.399	19.45
	PS	0.639	24.12	0.410	19.53
	LX	0.551	21.29	0.403	19.77
	SF	0.403	18.10	0.418	18.10
	All	0.728	26.41	0.440	19.98

It should be noted that we try our best to use different pre-training tasks from other techniques to ensure comprehensive comparisons, but it is not easy to adapt all of the pretext tasks for all of the encoders, e.g., MNP in Code2vec/Code2seq, Mask LM in Transformer. It is because MNP is designed specifically for Code2vec/Code2seq. Adapting TBCNN/Bi-LSTM for MNP-pretraining depends on various factors and requires certain configurations such as choosing an AST parser, processing paths in ASTs, choosing the right parameters for the prediction layer, etc. Note that the choice of AST parsers alone could affect the performance of programming-language processing tasks significantly [14]. There is no guarantee of the best settings for completely unbiased comparisons if we adapt. The same reasons are applicable to Masked-LM, which is designed specifically for the Transformer to process sequences. Thus, we only chose to adapt the pre-training task designed specifically for an encoder (MNP for Code2vec, Masked-LM for Transformer).

## 6 Analysis and Ablation Study

We perform some analysis and ablation studies to measure how different design choices can affect the performance of Corder.

### 6.1 Impact of Transformation Operators

We carry out an ablation study to evaluate how each transformation operator affects the performance of particular code learning tasks. We perform separate training of our Corder pretext task using different transformation operators in our Corder training algorithm. In particular, when taking the transformation operators to transform the code snippet, the set of available operators  $\mathcal{T}$  only contain one single operator.

An issue with the comparison is that the number of code snippets per operator may vary, leading to unfair comparisons between the

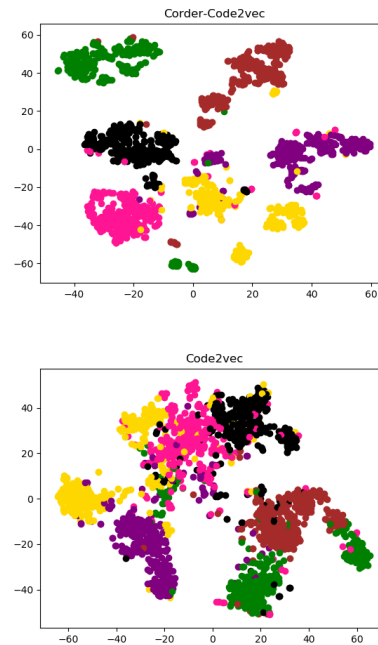
operators. This is because an operator is applied to modify the syntax and semantic of the code based on certain constraints, but such constraints may not apply all the time in the code. For example, the SF requires the snippet to contain at least one switch statement, but not all the snippets contain the switch statement. On the other hand, most of the snippets are applicable for VR because it is easy to change the variable names. Concretely, for JavaMed, the number of snippets applicable for VR, PS, US, LX, SF, are 919823, 64953, 352443, 26042, 3213, respectively. To see the impact of the operators, we perform this analysis under two settings: (1) we perform the transformation for each of the operators on the original number of snippets; and (2) we only select 3213 snippets for VR, PS, and US, in which the snippets must be the same as the snippets applicable for SF. We train Corder on the encoders with similar settings in the Evaluation Section, but we only use one operator at a time. Then we perform the fine-tuning process on the text-to-code retrieval and code summarization task, also similar to the Evaluation Section. The results of this analysis can be seen in Table 5. The column "Original" means the analysis is performed when using all of the snippets. The columns "Downsampled" means the analysis is performed when downsampling the number of snippets for VR, US, PS into the same number of snippets of SF. There are a few observations:

- In the "Original" setting, although VR has the most number of applicable snippets, its performance is among the worst, for either Text-to-Code Retrieval or Code Summarization. PS and US are two operators that perform the best (PS is slightly better most of the time). The performance of SF is comparable to VR but the number of snippets for SF is much fewer.
- In the "Downsampled" setting, SF becomes as comparable to PS and US, and these 3 operators perform much better than the VR, either in Text-to-Code Retrieval or Code Summarization. SF also performs the best in Text-to-Code Retrieval.

With the observations, we can conclude that changing the code structures is crucial to learn a decent model of source code. VR only changes the text information, while PS, US, and SF modify the code structure extensively.

## 6.2 Visualization for Code Embeddings

We visualize the code vectors to help understand and explain why the vectors produced by Corder pre-training are better than the vectors produced by other We choose Code2vec as the encoder for this analysis since Code2vec has been adapted in two different pretext tasks: (1) Corder pretext task (Corder-Code2vec); and (2) method name prediction task [6] (Code2vec). The goal is to show that our Corder pretext task performs better than the method name prediction as a pretext task to train the source code model. We use the code snippets in OJ dataset [55] that has been used for the code-to-code retrieval task. We randomly choose the embeddings of the first 6 classes of the OJ dataset then we use T-SNE [51] to reduce the dimensionality of the vectors into two-dimensional space and visualize. As shown in Figure 5, the vectors produced by Corder-Code2vec group similar code snippets into the same cluster with much clearer boundaries. This means that our instance discrimination task is a better pretext task than the method name prediction task in Alon et al. [6] for the same Code2vec encoder.



**Figure 5: Visualization of the vector representations of the code snippets from 6 classes in the OJ Dataset produced by Corder-Code2vec and MNP-Code2vec**

## 7 Conclusion

We have proposed Corder, a self-supervised learning approach that can leverage large-scale unlabeled data of source code. Corder works by training the neural network over a contrastive learning objective to compare similar and dissimilar code snippets that are generated from a set of program transformation operators. The snippets produced by such operators are syntactically diverse but semantically equivalent. The goal of the contrastive learning method is to minimize the distance between the representations of similar snippets (positives) and maximize the distance between dissimilar snippets (negatives). We have adapted Corder for three tasks: code-to-code retrieval, fine-tuning for text-to-code retrieval, fine-tuning for code summarization, and have found that Corder pre-training significantly outperforms other models not using contrastive learning on these three tasks.

## Acknowledgements

This research is supported by the Singapore Ministry of Education (MOE) Academic Research Fund (AcRF) Tier 2 Award No. MOE2019-T2-1-193 and RISE Lab Operational Fund from SCIS at SMU, Royal Society projects (IES/R1/191138, IES/R3/193175), EPSRC STRIDE project (EP/T017465/1), and Huawei Trustworthy Software Engineering and Open Source Lab. We also thank the anonymous reviewers for their insightful comments and suggestions, and thank the authors of related work for sharing data.

## References

- [1] [n.d.]. Krugle Code Search. <https://krugle.com/>. Accessed: 2020-09-30.

- [2] Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2020. A transformer-based approach for source code summarization. *arXiv preprint arXiv:2005.00653* (2020).
- [3] Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. 2018. Learning to Represent Programs with Graphs. In *ICLR*.
- [4] Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav. 2019. code2seq: Generating Sequences from Structured Representations of Code. In *ICLR*.
- [5] Uri Alon, Roy Sadaka, Omer Levy, and Eran Yahav. 2020. Structural language models of code. In *International Conference on Machine Learning*. PMLR, 245–256.
- [6] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. Code2Vec: Learning Distributed Representations of Code. In *POPL*. 40:1–40:29.
- [7] Antonio Valerio Miceli Barone and Rico Sennrich. 2017. A Parallel Corpus of Python Functions and Documentation Strings for Automated Code Documentation and Code Generation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017, Volume 2: Short Papers*, Greg Kondrak and Taro Watanabe (Eds.). Asian Federation of Natural Language Processing, 314–319. <https://www.aclweb.org/anthology/I17-2053/>
- [8] Marc Brockschmidt, Miltiadis Allamanis, Alexander L. Gaunt, and Oleksandr Polozov. 2019. Generative Code Modeling with Graphs. In *7th ICLR*.
- [9] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. 1994. Signature verification using a " siamese" time delay neural network. In *Advances in neural information processing systems*. 737–744.
- [10] Nghi DQ Bui, Lingxiao Jiang, and Yijun Yu. 2017. Cross-language learning for program classification using bilateral tree-based convolutional neural networks. *arXiv preprint arXiv:1710.06159* (2017).
- [11] Nghi DQ Bui, Yijun Yu, and Lingxiao Jiang. 2019. Bilateral dependency neural networks for cross-language algorithm classification. In *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 422–433.
- [12] Nghi DQ Bui, Yijun Yu, and Lingxiao Jiang. 2019. SAR: learning cross-language API mappings with little knowledge. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 796–806.
- [13] Nghi DQ Bui, Yijun Yu, and Lingxiao Jiang. 2020. InferCode: Self-Supervised Learning of Code Representations by Predicting Subtrees. *arXiv e-prints* (2020), arXiv–2012.
- [14] Nghi DQ Bui, Yijun Yu, and Lingxiao Jiang. 2020. TreeCaps: Tree-Based Capsule Networks for Source Code Processing. *arXiv preprint arXiv:2009.09777* (2020).
- [15] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709* (2020).
- [16] Xinyun Chen, Chang Liu, and Dawn Song. 2018. Tree-to-tree neural networks for program translation. In *NeurIPS*. 2547–2557.
- [17] Zimin Chen and Martin Monperrus. 2019. A literature study of embeddings on source code. *arXiv preprint arXiv:1904.03061* (2019).
- [18] Xiao Cheng, Haoyu Wang, Jiayi Hua, Guoai Xu, and Yulei Sui. 2021. DeepWukong: Statically detecting software vulnerabilities using deep graph neural network. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 30, 3 (2021), 1–33.
- [19] Michael L Collard, Michael John Decker, and Jonathan I Maletic. 2013. srcml: An infrastructure for the exploration, analysis, and manipulation of source code: A tool demonstration. In *ICSM*. 516–519.
- [20] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2017. Supervised learning of universal sentence representations from natural language inference data. *arXiv preprint arXiv:1705.02364* (2017).
- [21] George E Dahl, Jack W Stokes, Li Deng, and Dong Yu. 2013. Large-scale malware classification using random projections and neural networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. 3422–3426.
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [23] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).
- [24] Patrick Fernandes, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Structured Neural Summarization. In *7th ICLR*.
- [25] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. 2017. Self-supervised video representation learning with odd-one-out networks. 3636–3645.
- [26] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. 2018. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728* (2018).
- [27] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep code search. In *40th ICSE*. 933–944.
- [28] Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. 2017. DeepAM: Migrate APIs with Multi-modal Sequence to Sequence Learning. In *IJCAI* (Melbourne, Australia). 3675–3681.
- [29] Rahul Gupta, Aditya Kanade, and Shirish Shevade. 2019. Neural Attribution for Semantic Bug-Localization in Student Programs. In *NeurIPS*. 11861–11871.
- [30] Jingxuan He, Cheng-Chun Lee, Veselin Raychev, and Martin Vechev. 2021. Learning to Find Naming Issues with Big Code and Small Supervision. (2021).
- [31] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
- [32] Xing Hu, Ge Li, Xin Xia, David Lo, and Zhi Jin. 2018. Deep code comment generation. 200–210.
- [33] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436* (2019).
- [34] Yasir Hussain, Zhiqiu Huang, Yu Zhou, and Senzhang Wang. 2020. Deep transfer learning for source code modeling. *International Journal of Software Engineering and Knowledge Engineering* 30, 05 (2020), 649–668.
- [35] Bill Ingram. 2018. *A Comparative Study of Various Code Embeddings in Software Semantic Matching*. <https://github.com/waigram/code-embeddings>.
- [36] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2016. Summarizing source code using a neural attention model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2073–2083.
- [37] Darryl Jarman, Jeffrey Berry, Riley Smith, Ferdian Thung, and David Lo. 2021. Legion: Massively Composing Rankers for Improved Bug Localization at Adobe. *IEEE Transactions on Software Engineering* (2021).
- [38] Lin Jiang, Hui Liu, and He Jiang. 2019. Machine learning based recommendation of method names: how far are we. In *34th ASE*. 602–614.
- [39] Hong Jin Kang, Tegawendé F Bissyandé, and David Lo. 2019. Assessing the generalizability of code2vec token embeddings. In *34th ASE*. 1–12.
- [40] Dahun Kim, Donghyeon Cho, and In So Kweon. 2019. Self-supervised video representation learning with space-time cubic puzzles. In *AAAI*, Vol. 33. 8545–8552.
- [41] Kisub Kim, Dongsun Kim, Tegawendé F Bissyandé, Eunjong Choi, Li Li, Jacques Klein, and Yves Le Traon. 2018. FaCoY: a code-to-code search engine. 946–957.
- [42] Ryan Kiros, Yunkun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. In *NeurIPS*. 3294–3302.
- [43] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M Rush. 2017. Opennmt: Open-source toolkit for neural machine translation. *arXiv preprint arXiv:1701.02810* (2017).
- [44] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*. Association for Computational Linguistics, 177–180.
- [45] Lingpeng Kong, Cyprien de Masson d’Autume, Wang Ling, Lei Yu, Zihang Dai, and Dani Yogatama. 2019. A mutual information maximization perspective of language representation learning. *arXiv preprint arXiv:1910.08350* (2019).
- [46] Bruno Korbar, Du Tran, and Lorenzo Torresani. 2018. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*. 7763–7774.
- [47] Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. 1188–1196.
- [48] Jian Li, Pinjia He, Jieming Zhu, and Michael R Lyu. 2017. Software defect prediction via convolutional neural network. In *IEEE QRS*. 318–328.
- [49] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2016. Gated Graph Sequence Neural Networks. In *ICLR*.
- [50] Lajanugen Logeswaran and Honglak Lee. 2018. An efficient framework for learning sentence representations. *arXiv preprint arXiv:1803.02893* (2018).
- [51] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of Machine Learning Research* 9, Nov (2008), 2579–2605.
- [52] Aravindh Mahendran, James Thewlis, and Andrea Vedaldi. 2018. Cross pixel optical-flow similarity for self-supervised learning. In *Asian Conference on Computer Vision*. 99–116.
- [53] Henry Massalin. 1987. Superoptimizer - A Look at the Smallest Program. In *Proceedings of the Second International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS II), Palo Alto, California, USA, October 5-8, 1987*, Randy H. Katz and Martin Freeman (Eds.). ACM Press, 122–126. <https://dl.acm.org/citation.cfm?id=36194>
- [54] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NeurIPS*. 3111–3119.
- [55] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional neural networks over tree structures for programming language processing. In *AAAI*.
- [56] R. Nix and J. Zhang. 2017. Classification of Android apps and malware using deep neural networks. In *International Joint Conference on Neural Networks*. 1871–1878.
- [57] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [58] Michael Pradel and Koushik Sen. 2018. DeepBugs: A learning approach to name-based bug detection. *ACM on Programming Languages* 2, OOPSLA (2018), 147.

- [59] Md Rabin, Rafiqul Islam, Nghi DQ Bui, Yijun Yu, Lingxiao Jiang, and Mohamad Amin Alipour. 2020. On the Generalizability of Neural Program Analyzers with respect to Semantic-Preserving Program Transformations. *arXiv preprint arXiv:2008.01566* (2020).
- [60] Saksham Sachdev, Hongyu Li, Sifei Luan, Seohyun Kim, Koushik Sen, and Satish Chandra. 2018. Retrieval on Source Code: A Neural Code Search. In *2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages* (Philadelphia, PA, USA), 31–41. <https://doi.org/10.1145/3211346.3211353>
- [61] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 815–823.
- [62] Tian Shi, Liuqing Li, Ping Wang, and Chandan K Reddy. 2020. A Simple and Effective Self-Supervised Contrastive Learning Framework for Aspect Detection. *arXiv preprint arXiv:2009.09107* (2020).
- [63] Yulei Sui, Xiao Cheng, Guanqin Zhang, and Haoyu Wang. 2020. Flow2Vec: value-flow-based precise code embedding. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 1–27.
- [64] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*. 3104–3112.
- [65] Jeffrey Svajlenko and Chanchal K Roy. 2015. Evaluating clone detection tools with bigclonebench. In *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 131–140.
- [66] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. 2020. Intellicode compose: Code generation using transformer. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1433–1443.
- [67] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [68] Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S. Yu. 2018. Improving Automatic Source Code Summarization via Deep Reinforcement Learning. In *33rd ASE* (Montpellier, France). New York, NY, USA, 397–407. <https://doi.org/10.1145/3238147.3238206>
- [69] Hong Wang, Xin Wang, Wenhan Xiong, Mo Yu, Xiaoxiao Guo, Shiyu Chang, and William Yang Wang. 2019. Self-supervised learning for contextualized extractive summarization. *arXiv preprint arXiv:1906.04466* (2019).
- [70] Haoye Wang, Xin Xia, David Lo, John Grundy, and Xinyu Wang. 2021. Automatic Solution Summarization for Crash Bugs. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 1286–1297.
- [71] Jiawei Wu, Xin Wang, and William Yang Wang. 2019. Self-supervised dialogue learning. *arXiv preprint arXiv:1907.00448* (2019).
- [72] Xinli Yang, David Lo, Xin Xia, Yun Zhang, and Jianling Sun. 2015. Deep Learning for Just-in-Time Defect Prediction. In *IEEE QRS*. 17–26.
- [73] Michihiro Yasunaga and Percy Liang. 2020. Graph-based, Self-Supervised Program Repair from Diagnostic Feedback. *arXiv preprint arXiv:2005.10636* (2020).
- [74] Huangzhao Zhang, Zhuo Li, Ge Li, Lei Ma, Yang Liu, and Zhi Jin. 2020. Generating Adversarial Examples for Holding Robustness of Source Code Processing Models. In *34th AAAI*.
- [75] Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. 2019. A novel neural source code representation based on abstract syntax tree. In *41st ICSE*. 783–794.
- [76] Richard Zhang, Phillip Isola, and Alexei A Efros. 2016. Colorful image colorization. In *European conference on computer vision*. 649–666.
- [77] Han Zhao, Zhengdong Lu, and Pascal Poupart. 2015. Self-adaptive hierarchical sentence model. *arXiv preprint arXiv:1504.05070* (2015).
- [78] Yaqin Zhou, Shangqing Liu, Jing Kai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective Vulnerability Identification by Learning Comprehensive Program Semantics via Graph Neural Networks. In *NeurIPS*. 10197–10207.