

CRaDLe: Deep code retrieval based on semantic Dependency Learning

Wenchao Gu^a, Zongjie Li^b, Cuiyun Gao^{b,*}, Chaozheng Wang^b, Hongyu Zhang^c,
Zenglin Xu^b, Michael R. Lyu^a

^a The Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, China

^b The School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, China

^c The University of Newcastle, Australia

ARTICLE INFO

Article history:

Received 10 November 2020

Received in revised form 8 March 2021

Accepted 15 April 2021

Available online 26 April 2021

Keywords:

Code retrieval

Semantic dependency

Dependency learning

Neural network

ABSTRACT

Code retrieval is a common practice for programmers to reuse existing code snippets in the open-source repositories. Given a user query (i.e., a natural language description), code retrieval aims at searching the most relevant ones from a set of code snippets. The main challenge of effective code retrieval lies in mitigating the semantic gap between natural language descriptions and code snippets. With the ever-increasing amount of available open-source code, recent studies resort to neural networks to learn the semantic matching relationships between the two sources. The statement-level dependency information, which highlights the dependency relations among the program statements during the execution, reflects the structural importance of one statement in the code, which is favorable for accurately capturing the code semantics but has never been explored for the code retrieval task. In this paper, we propose CRaDLe, a novel approach for Code Retrieval based on statement-level semantic Dependency Learning. Specifically, CRaDLe distills code representations through fusing both the dependency and semantic information at the statement level, and then learns a unified vector representation for each code and description pair for modeling the matching relationship. Comprehensive experiments and analysis on real-world datasets show that the proposed approach can accurately retrieve code snippets for a given query and significantly outperform the state-of-the-art approaches on the task.

© 2021 Elsevier Ltd. All rights reserved.

1. Introduction

Implementing projects from scratch is tedious for programmers. In most cases, they know what they want to do, but do not have the capability to implement all the details. For example, a Python programmer may want to “convert `date_string` into `datetime format`”, but not able to recognize the proper syntax `datetime.strptime(date_string, format)` for the realization. To mitigate the impasse, it is common for programmers to search the web in natural language (NL), find relevant code snippets, and modify them into the desired form (Brandt, Guo, Lewenstein, Dontcheva, & Klemmer, 2009). Many code retrieval approaches (Brandt et al., 2009; Lv et al., 2015; McMillan, Grechanik, Poshyanyk, Xie, & Fu, 2011) have been proposed to improve the recommendation accuracy of the returned code snippets given a natural language description. The main challenge of effective code retrieval is the semantic gap between source

code and natural language descriptions since the two sources are heterogeneous and share few common lexical tokens, synonyms and language structures (Gu, Zhang, & Kim, 2018).

Prior efforts have been conducted for effective code retrieval. The existing research can be divided into two categories according to the involved techniques, i.e., Information Retrieval (IR)-based and Deep Neural Network (DNN)-based. The IR-based techniques rely on token-wise similarities between source code and queries. Since the variable and API definitions in code are generally word combinations or abbreviations in natural language, more semantically-similar tokens in code and queries can indicate more relevancy between them. For example, McMillan et al. propose Portfolio which utilizes keyword matching and PageRank to return a list of functions (McMillan et al., 2011). Lv et al. propose CodeHow to combine API matching for code retrieval (Lv et al., 2015). With an increasing amount of available source code and flourish development of deep learning techniques, many studies (Gu et al., 2018; Husain, Wu, Gazit, Allamanis, & Brockschmidt, 2019) propose to adopt neural network models for jointly encoding tokens of source code and queries into a single and joint vector space, where one encoder is employed for each input (natural or programming) sequence. The

* Corresponding author.

E-mail addresses: wcpu@cse.cuhk.edu.hk (W. Gu), lizongjie@stu.hit.edu.cn (Z. Li), gaocuiyun@hit.edu.cn (C. Gao), wangchaozheng@stu.hit.edu.cn (C. Wang), hongyu.zhang@newcastle.edu.au (H. Zhang), xuzenglin@hit.edu.cn (Z. Xu), lyu@cse.cuhk.edu.hk (M.R. Lyu).

objective is to map semantically relevant code and language into vectors that are near to each other in the vector space.

Considering the highly-structured characteristic of source code, recent research proposes to integrate the structural information of code such as Abstract Syntax Tree (AST) and Control Flow Graph (CFG) for representing code semantics (Wan et al., 2019; Yin & Neubig, 2017; Zhang et al., 2019), demonstrating the effectiveness of involving structural information for the task. However, the deep nature of the extracted trees in ASTs renders it hard for deep learning models to comprehensively capture the structural information (Zhang et al., 2019). CFG, which represents all possible execution paths for a program, may contain statement orders which are not contributing to the actual execution result, probably leading to biased code representation learning (Wan et al., 2019). In this paper, we propose to utilize statement-level dependency relations in a code snippet based on Program Dependency Graph (PDG). The PDG is established based on AST but less deeper than AST in the structure and only retains the execution paths that will affect the execution result. The dependency relations are then explicitly integrated with the statement-level semantics to capture the code semantics. Actually, the effectiveness of incorporating dependency relations for code representation learning has proven in tasks such as bug detection (Li, Wang, Nguyen, & Nguyen, 2019) and code clone detection (Henderson & Podgurski, 2016); while no prior work has explored the impact on the code retrieval task so far.

Specifically, we introduce a novel neural network model named CRaDLe, an abbreviation of **C**ode **R**etrieval based on **S**emantic **D**ependency **L**earning. CRaDLe couples both structural and semantic information of code at the statement level, where the code structures are extracted based on PDG. Extensive experiments have been conducted to verify the performance of the proposed approach. The evaluation results show that CRaDLe can significantly outperform the state-of-the-art models by at least 36.38% and 22.34% on two real datasets respectively, in terms of R@1, one standard metric for validating recommendation performance.

In summary, the main contributions of the paper include:

- We propose a novel code retrieval model, CRaDLe, to encode both source code and natural language queries into unified vector representations. CRaDLe is the first code retrieval approach that integrates the dependency and semantic information at the statement level for learning code representations.
- We conduct large-scale experimental evaluations on public benchmarks. The results demonstrate the superior performance of CRaDLe over the state-of-the-art and baseline models.

The rest of this paper is organized as follows. Section 2 introduces an overview of the proposed approach and details the design of the approach. Section 3 illustrates the experimental datasets, evaluation metrics, and implementation details. Section 4 elaborates on the experimental results. Section 6 surveys the related work and Section 7 concludes our work.

2. The proposed CRaDLe

In this section, we elaborate on the overview and detailed design of the proposed approach CRaDLe, including the code encoder, description encoder and the similarity measurement component.

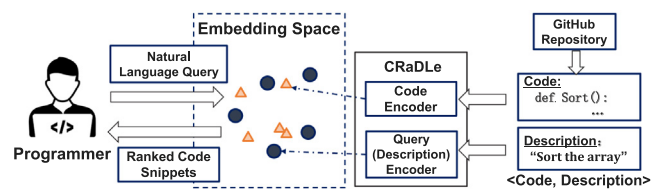


Fig. 1. Overview of the proposed CRaDLe.

2.1. Overview

Fig. 1 depicts the overview of the proposed approach, CRaDLe. The implementation includes both offline and online modes. During the offline stage, we first collect datasets containing `<code, description>` pairs. The collected code and descriptions are then preprocessed and separately encoded into vectors by the code encoder and query encoder respectively. Unified representations of code and corresponding descriptions are finally learnt after the offline training process, where semantically similar code and descriptions locate closely to each other in the same embedding space. During the online process, when a new natural language query arrives, the trained model recommends the most related code snippets to the programmer according to the semantic distances between code and the query in the embedding space.

Fig. 2 illustrates the overall framework of the CRaDLe approach, which details the design of the code encoder and description encoder. The code encoder fuses the statement-level semantics and distilled dependency information to represent the code semantics. The description encoder also embeds the token sequences in the descriptions to vectors. Finally, similarity matching scores between the code and descriptions are learnt based on their respective vector representations.

2.2. Code encoder

The code encoder aims at embedding code snippets into vector representations. We propose to integrate the statement-level token semantics with the dependency information between statements for accurately capturing the code semantics. We first illustrate the process conducted for the dependency information extraction, and then describe the networks proposed for learning statement-level dependency and semantic representations.

Algorithm 1 shows the procedures for the code encoder. The input of the code encoder includes the token matrix E comprised by a sequence of token embedding vectors $\{\mathbf{e}_{1,1}, \dots, \mathbf{e}_{i,j}, \dots\}$ and dependency matrix γ . First, the dependency embedding layer encodes the dependency matrix γ into dependency embeddings P . The token embedding layer then represents the token matrix E into statement-level representations T comprised by $\{\mathbf{t}_1, \dots, \mathbf{t}_i, \dots\}$. Finally, the token embeddings are concatenated with the dependency embeddings in statement level, and the newly comprised vectors are fed into the Bi-LSTM layer. The last hidden state vector from the Bi-LSTM layer is treated as the representation vector of the code.

2.2.1. Dependency information extraction

We obtain the dependency information between statements by adopting PDGs of the code snippets. PDG explicitly indicates the data dependency and control dependency of a program, where the data dependency can represent the relevant data flow relationships and control dependency exhibits the essential control flow relationships (Ferrante, Ottenstein, & Warren, 1987). Since there exists no mature tool for extracting PDG of one code

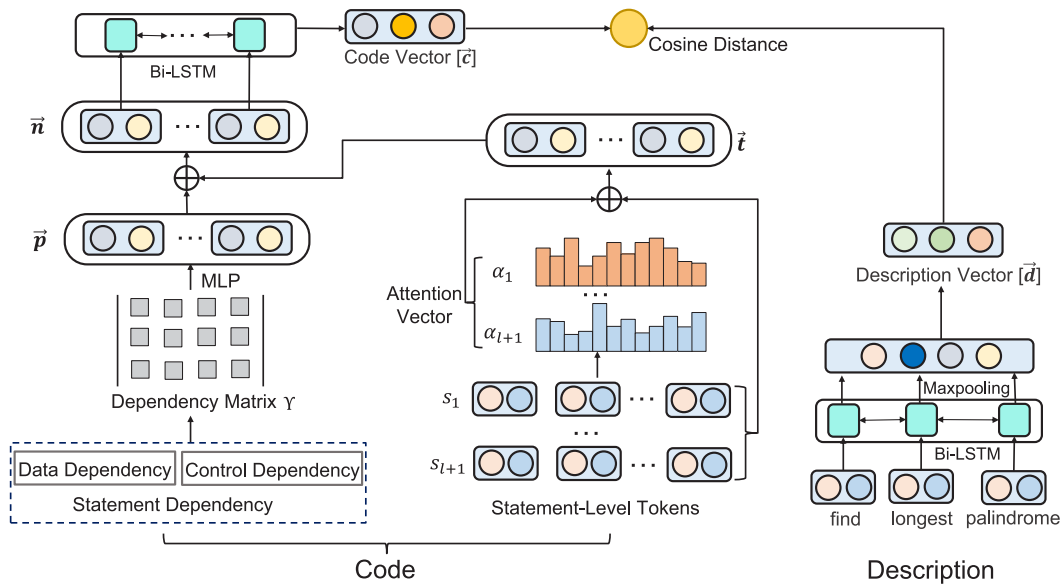


Fig. 2. Overall framework of the proposed CRaDLe.

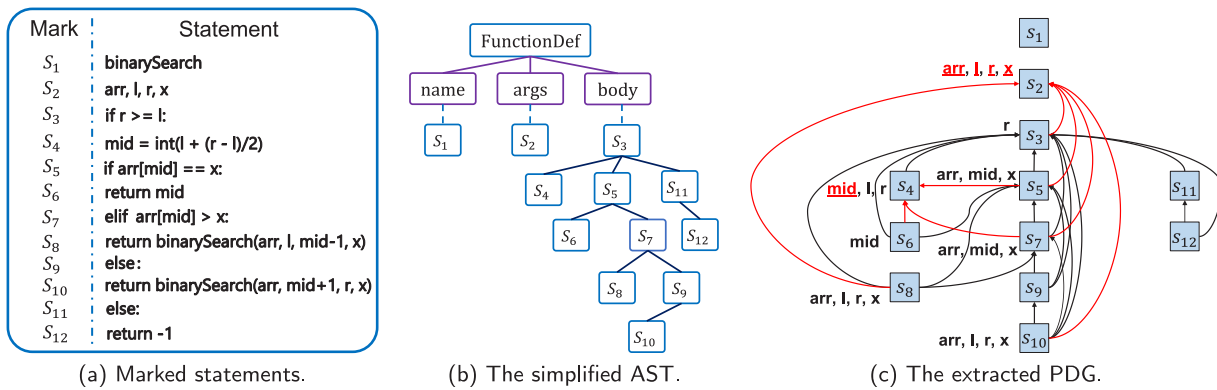


Fig. 3. Workflow for extracting PDG of the code snippet in Listing 1. For the extracted PDG in (c), red and black arrowed lines indicate data dependency and control dependency respectively. The tokens beside each statement block denote the variables (re)defined (highlighted in red underlined font) or used in the corresponding statement.

Algorithm 1: The algorithm of code encoding

input : the token matrix E , the matrix of input dependency: \mathcal{Y}
output: The representation vector of code: C
Function CODEENCODER(E, \mathcal{Y}):
 $P \leftarrow$ DependencyEmbedding(\mathcal{Y}); // corresponding to Eq. (1)
 $T \leftarrow$ TokenEmbedding(E);
 $S \leftarrow$ StatementAttention(T); // corresponding to Eqs. (2) and (3)
 $C \leftarrow$ SemanticDependencyEmbedding($[S; P]$);
 // corresponding to Eq. (4)
return C ;

snippet in interpreted languages such as Python, we propose to establish the PDG based on the AST of a code snippet.

For clarifying the PDG establishment process, we use the code example illustrated in Listing 1. Fig. 3(a) depicts the mark for each statement in the code example, in which we regard the function name and required parameters as two separate statements.

Function name can be treated as a short summary of the code functionality; while the definitions of the required parameters generally reflect the semantics of the input data. Treating function names and parameters separately could be helpful for capturing their respective semantics.

Fig. 3(b) demonstrates the simplified AST of the code example where we construct the AST in statement level and hide the details of each statement. The data dependency of one statement with the other statement can be identified if the variable used in one statement is (re)defined in the other statement and the value of the variable is unchanged on the execution path between these two statements. The control dependency of one statement with the other is determined if the execution of the statement relies on the execution results of the other one. The control dependency can be directly captured by the tree structure in the AST, i.e., statements in child leaf nodes are considered possessing dependent relations with the statements in the parent nodes. The extracted PDG is depicted in Fig. 3(c), with red arrowed lines and black arrowed lines indicate data dependency and control dependency between the two statements, respectively. Tokens beside each statement block denote the related variables, in which we use black or

red underlined variables to distinguish whether the variables are used or (re)defined in the corresponding statement. For example, the parent nodes of S_{10} in the AST include $S_3, S_5, S_7,$ and S_9 (as shown in Fig. 3(b)), so the control dependency between S_{10} and $S_{3,5,7,9}$ is marked in the obtained PDG. Also, the variable mid in S_6 , corresponding to line 5 in the Listing 1, is from S_4 , i.e., line 3 in the code example; so S_6 shows a data dependency relation to S_4 in the PDG.

2.2.2. Statement-level dependency embedding

The dependency embedding network is designed to encode the data dependency and control dependency involved in the PDG of a code snippet into a vector representation. According to the extracted PDG (as shown in Fig. 3(c)), we can build a dependency matrix $\mathcal{Y} \in \{0, 1\}^{(l) \times (l)}$, where l indicates the number of statements in the code. The element $v_{ij} = 1$ if the i th statement has a data/control dependency on the j th statement; otherwise $v_{ij} = 0$. Note that $v_{ij} \neq v_{ji}$. For example, S_4 and S_6 exhibit a data dependency relation, so $v_{64} = 1$. To embed the obtained dependency matrix \mathcal{Y} , we employ one layer of multi-layer perceptron (MLP):

$$\begin{aligned} \mathbf{p}_i &= \tanh(\mathbf{W}^T \mathbf{v}_i), \forall i = 1, 2, \dots, l, \\ \mathbf{P} &= [\mathbf{p}_1, \dots, \mathbf{p}_l], \end{aligned} \quad (1)$$

where \mathbf{W}^T is the matrix of trainable parameters in MLP and \mathbf{p}_i is the embedding of the dependency information for each statement.

2.2.3. Statement-level token embedding

The token embedding network is designed for capturing the semantics of each statement based on the constituted tokens. We first tokenize the statements into sequences of tokens following Gu et al.'s work (Gu et al., 2018), during which process duplicate tokens and the keywords in the programming language such as `while` and `break` are removed. Then tokens in each sequence are embedded into vectors individually through an embedding layer. An attention layer is utilized to compute a weighted average. Given a sequence of token embedding vectors $\{\mathbf{e}_{i,1}, \dots, \mathbf{e}_{i,j}, \dots\}$ for the i th statement, the attention weight $\alpha_{i,j}$ for each $\mathbf{e}_{i,j}$ is calculated as follows:

$$\alpha_{i,j} = \frac{\exp(\mathbf{e}_{i,j}^T)}{\sum_j \exp(\mathbf{e}_{i,j}^T)}. \quad (2)$$

Each statement is embedded based on the attention weights $\alpha_{i,j}$.

$$\mathbf{t}_i = \sum_j \alpha_{i,j} \mathbf{e}_{i,j}^T, \quad (3)$$

where i indicates the i th statement.

2.2.4. Semantic dependency embedding

We consider both statement-level dependency and semantic information for learning the vector representation of a code snippet. Specifically, for each statement s_i , we concatenate its dependency embedding \mathbf{p}_i and token embedding \mathbf{t}_i as the representation of the statement, i.e., $\mathbf{s}_i = [\mathbf{t}_i; \mathbf{p}_i]$. We finally adopt bi-LSTM to encode the sequence of the statement embeddings and use the last hidden state as the vector representation of the code.

$$\mathbf{c} = \text{BiLSTM}(\mathbf{h}_l, \mathbf{s}_l), \quad (4)$$

where l indicates the number of statements.

```

1 def binarySearch(arr, l, r, x):
2     if r >= l:
3         mid = int((l + (r - 1))/2)
4         if arr[mid] == x:
5             return mid
6         elif arr[mid] > x:
7             return binarySearch(arr, l, mid-1, x)
8         else:
9             return binarySearch(arr, mid+1, r, x)
10    else:
11        return -1

```

Code Listing 1: An example of Python code snippet for illustrating the semantic dependency learning process.

2.3. Description encoder

The description encoder aims at embedding natural language descriptions into vectors. Given a description $D = \{w_1, \dots, w_k, \dots, w_{N_d}\}$ comprising a sequence of N_d words, the description encoder embeds it into a vector \mathbf{d} using a bi-LSTM model with maxpooling:

$$\begin{aligned} \mathbf{h}_k &= \text{BiLSTM}(\mathbf{h}_{k-1}, \mathbf{w}_k), \forall k = 1, 2, \dots, N_d, \\ \mathbf{d} &= \text{maxpooling}([\mathbf{h}_1, \dots, \mathbf{h}_{N_d}]). \end{aligned} \quad (5)$$

The maxpooling layer is used to mitigate the effect of long-term information loss caused by the LSTM mechanism and catch the global feature of the whole sentence.

2.4. Similarity measurement

The semantic similarity between the code vector \mathbf{c} and description vector \mathbf{d} is calculated based on its cosine distance in the embedding space:

$$\cos(\mathbf{c}, \mathbf{d}) = \frac{\mathbf{c}^T \mathbf{d}}{\|\mathbf{c}\| \|\mathbf{d}\|}. \quad (6)$$

The vector features of the two different embedding models are trained using the loss function, i.e., Eq. (6), to maximize the cosine similarities in the projected space, so aligned code and descriptions would be close to each other in the space. Such design is widely adopted in prior code search studies (Cambronero, Li, Kim, Sen, & Chandra, 2019; Gu et al., 2018; Sachdev et al., 2018). The target of the design is to get unified representations for both code and description, so as to mitigate the problem of semantic gap between them. The higher the similarity, the more relevant the code is to the description.

2.5. Model training

We obtain the representation vectors for code snippets and descriptions based on the proposed code encoder and description encoder, respectively. Following previous studies (Cambronero et al., 2019; Gu et al., 2018; Sachdev et al., 2018), we project the code vectors and description vectors to the same space, and train the vectors for aligned code snippets and descriptions to be close in the space.

Specifically, every single code snippet in the training data T will be constructed as a triplet $\langle C, D+, D- \rangle$. C represents the code snippet from the training Corpora, $D+$ indicates the description which semantically matches the code snippet in the ground truth, and $D-$ denotes the negative description which is randomly chosen from the training corpora with the true description excluded. The loss function is as below:

$$\mathcal{L}(\theta) = \sum_{(C, D+, D-) \in T} \max(0, \epsilon - \cos(\mathbf{c}, \mathbf{d}+) + \cos(\mathbf{c}, \mathbf{d}-)), \quad (7)$$

where θ denotes the parameters in the proposed model, \mathbf{c} denotes the code vector of C , $\mathbf{d}+$ and $\mathbf{d}-$ denote the description vectors

Table 1

Statistics of the number of statements in CodeSearchNet dataset.

#Statements	Training set	Validation set	Test set
0 ~ 10	2,30,183	12,413	12,326
11 ~ 20	1,17,060	6,364	6,361
21 ~ 30	32,904	1,875	1,843
31 ~ 40	12,834	755	633
41 ~ 50	5,723	386	326
51 ~	8,422	509	413

Table 2

Statistics of the number of statements in Code2Seq dataset.

#Statements	Training set	Validation set	Test set
0 ~ 10	2,18,679	32,429	33,210
11 ~ 20	73,870	11,301	12,251
21 ~ 30	20,956	3,215	3,478
31 ~ 40	7,957	1,251	1,370
41 ~ 50	3,540	573	632
51 ~	4,326	650	786

of $D+$ and $D-$, respectively. Based on the training loss function, we can get unified representations for both code and description, thus mitigating the semantic gap between them.

3. Experimental setup

In the section, we introduce the collected dataset for experimentation, the evaluation metrics, implementation details and baseline models.

3.1. Dataset collection

Two datasets are adopted for our experimental evaluation. One dataset is obtained from CodeSearchNet (Husain et al., 2019), a publicly-available GitHub repository. We focus on the Python program language since it is one of the most popular programming languages, accounting for more than 30% of the total market share as PYPL reported (PYPL, 2020). Detailed statistics of the dataset can be found in Table 3. All the code in the corpus is in Python and with English descriptions. We have 407,126, 22,302, and 21,902 <code, description> pairs for training, validating and testing, respectively. The median and average numbers of the statements in the code are around 10. We also observe that the statements contain around three tokens on average, with the minimum at zero which is because the input parameters beside the method name are treated as an individual statement and some code snippets may not require any input parameters. Another dataset is from Code2seq (Alon, Brody, Levy, & Yahav, 2019), with the statistics illustrated in Table 4. We only select the code written in Python 3 from both datasets since the PDG extraction tool (introduced in Section 2.2.1) is specifically designed for Python 3 and may fail to parse the code written in Python 2.

Tables 1 and 2 illustrate the distribution of statements numbers of the codes in the two dataset, i.e., CodeSearchNet and Code2Seq, respectively. We can observe that the long tail phenomenon occurs in the two datasets. Besides, more than 50% of the code has ≤ 10 statements and more than 80% has ≤ 20 statements.

3.2. Performance measurement

Following the evaluation settings in Wan et al. (2019), we fix a set of 999 distractor snippets \mathbf{c}_j for each test pair $(\mathbf{c}_i, \mathbf{d}_i)$ and calculate the average ranking score for all the testing pairs as the evaluation result. We involve two metrics: $R@k$ and MRR, for validating the ranking performance.

Table 3

Statistics of the CodeSearchNet dataset.

# <code, description>	Training	Validating	Testing
	4,07,126	22,302	21,902
Statistics of # statements in code			
Min.	1	1	1
Med.	7	8	7
Max.	1,385	909	363
Ave.	11.45	11.87	11.24
Statistics of # tokens in the statements			
Min.	0	0	0
Med.	3	3	3
Max.	514	155	83
Ave.	3.92	3.87	3.91

Table 4

Statistics of the Code2seq dataset.

# <code, description>	Training	Validating	Testing
	3,29,328	49,419	51,727
Statistics of # statements in code			
Min.	2	2	2
Med.	7	7	7
Max.	1,463	416	1,463
Ave.	10.17	10.33	10.68
Statistics of # tokens in the statements			
Min.	0	0	0
Med.	3	3	3
Max.	682	199	1,864
Ave.	3.75	3.73	3.73

3.2.1. $R@k$

$R@k$ is a common metric to evaluate whether an approach can retrieve the correct answer in the top k returning results. It is widely used by many studies on the code retrieval task. The metric is calculated as follows:

$$R@k = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \delta(\text{FRank}_q < k), \quad (8)$$

where Q denotes the query set and FRank_q denotes the rank of the correct answer for query q . The function $\delta(\text{FRank}_q < k)$ returns 1 if the rank of the correct answer within the top k returning results otherwise the function returns 0. A higher $R@k$ indicates a better code retrieval performance.

3.2.2. MRR

Mean Reciprocal Rank (MRR) is the average of the reciprocal ranks of the correct answers of query set Q , which is another popular evaluation metric for the code retrieval task. The metric MRR is calculated as follows:

$$\text{MRR} = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \frac{1}{\text{FRank}_q}. \quad (9)$$

The higher the MRR value is, the better performance the model has.

3.3. Implementation details

In our experiment, we select the top 10,000 words according to the word frequencies as the vocabularies of code snippets and descriptions, respectively. All the word embeddings are randomly initialized and adjusted during training. The dimension of word embedding is set as 256. All LSTMs have 1024 hidden units in each direction. The maximum number of considered statements

in the code and the maximum number of tokens in each statement are set as 20 and 5, respectively. The sequence lengths of descriptions are limited as 30 following the work (Gu et al., 2018). The CRaDLe model is trained via the AdamW algorithm (Kingma & Ba, 2015) and the learning rate is $2.08e-4$. To mitigate the overfitting issue, we add a dropout layer with dropout rate at 0.25. We train our models on a server with one Nvidia GeForce RTX 2080 Ti and 11 GB memory. The training lasts ~ 20 h with 200 epochs and the early stopping strategy (Goodfellow, Bengio, Courville, & Bengio, 2016) is adopted to avoid overfitting.

3.4. Baseline models

We compare our proposed model with several state-of-the-art baseline models. **CODEnn** is one of the state-of-the-art models proposed in Gu et al. (2018). This model extracts the method name, API sequence and tokens from the code and utilizes neural network to learn the unified vector representation of query and these code features. **UNIF** (Cambroner et al., 2019) focuses on the semantic information from the tokens in the code and utilizes embedding techniques and attention mechanism to embed the tokens in the query and code into a single vector respectively. The projection of the query and code vector in the same space is learned by this model. **NeuralBoW** (Wang & Manning, 2012) embeds each token in the two input sequences to a learnable embedding. The token embeddings are then combined into a sequence embedding using max-pooling and an attention-like weighted sum mechanism. The **RNN** baseline adopts two-layer bi-directional LSTM model (Cho, van Merriënboer, Bahdanau, & Bengio, 2014) to encode the input sequences. **CONV** (Kim, 2014) uses 1D convolutional neural network over both the input sequences of tokens. **CONVSelf** (Lin et al., 2017) combines 1D convolutional neural network and self-attention layer to embed both input sequences. **SelfAttn** (Husain et al., 2019) utilizes multi-head attention (Vaswani et al., 2017) to encode both input sequences of tokens, and has proven effective on multiple types of programming languages such as Python and JavaScript. The hyperparameters of the baselines are defined according to the original papers (Cambroner et al., 2019; Gu et al., 2018; Husain et al., 2019). During implementing CODEnn, NeuralBoW, RNN, CONV, CONVSelf and SelfAttn, we directly utilized the released code; while for UNIF, we tried our best to replicate the code according to the paper and will make the replication publicly available.

4. Experimental results

In this section, we present the evaluation results, including the main results, parameter analysis, case studies and error analysis.

4.1. Main results

Involving semantic dependency embeddings increases the code search performance. Tables 5 and 6 illustrate the evaluation results comparing with the baseline models. As can be seen, CRaDLe presents the best performance comparing with all the baseline models, increasing the performance of 36.38% in terms of $R@1$, 17.13% in terms of $R@5$, 12.54% in terms of $R@10$ and 25.26% in terms of MRR at least on the dataset of CodeSearchNet. CRaDLe can achieve the improvement of the performance at least 22.34%, 22.51%, 21.54% and 21.79% in $R@1$, $R@5$, $R@10$ and MRR on the dataset of Code2Seq, respectively. This indicates that CRaDLe can rank the correct answer the top more accurately when given a natural language query. The improvement on $R@1$ is most significant among all the metrics in our proposed model, which is over 20% in both datasets. $R@1$ is the metric concerned most by programmers since they prefer to use the code search

Table 5

Comparison results with baseline models on the CodeSearchNet dataset. The best results are highlighted in **bold** fonts.

Approach	R@1	R@5	R@10	MRR
CODEnn	0.367	0.573	0.652	0.465
UNIF	0.379	0.615	0.706	0.490
NeuralBoW	0.521	0.747	0.807	0.622
RNN	0.556	0.772	0.832	0.654
CONV	0.475	0.703	0.776	0.579
CONVSelf	0.571	0.788	0.845	0.668
SelfAttn	0.580	0.786	0.840	0.673
CRaDLe _{maxpooling}	0.777	0.914	0.946	0.838
CRaDLe	0.791	0.923	0.951	0.843

Table 6

Comparison results with baseline models on the Code2seq dataset. The best results are highlighted in **bold** fonts.

Approach	R@1	R@5	R@10	MRR
CODEnn	0.330	0.532	0.617	0.427
UNIF	0.380	0.588	0.668	0.478
NeuralBoW	0.546	0.693	0.738	0.615
RNN	0.438	0.623	0.688	0.526
CONV	0.425	0.584	0.645	0.502
CONVSelf	0.470	0.642	0.700	0.552
SelfAttn	0.525	0.683	0.731	0.599
CRaDLe _{maxpooling}	0.664	0.843	0.892	0.745
CRaDLe	0.668	0.849	0.897	0.749

system which can return the best results in first. The higher MRR score further verifies the effectiveness of CRaDLe. The difference between CRaDLe and the baseline models is the code representation strategy, which shows the effectiveness of the semantic dependency embeddings for code search.

Attention mechanism can be helpful for effective code search. By comparing CONV with CONVSelf, we can observe that with the attention mechanism integrated, CONV presents a better performance than the pure CONV model on both datasets. For example, CONVSelf increases the accuracy of CONV by 20.21% and 15.37% in terms of $R@1$ and MRR on the CodeSearchNet dataset, respectively. Similar result also appears on the Code2Seq dataset. The results imply the effectiveness of the attention mechanism on the code search task. We also compared with the performance of the CRaDLe_{maxpooling} where the attention mechanism is replaced with the max pooling strategy (Lee, Gallagher, & Tu, 2016). As can be seen in Tables 5 and 6, CRaDLe with the attention mechanism involved outperforms the CRaDLe with max pooling strategy integrated on both datasets, which further demonstrates the effectiveness of the attention mechanism on the task.

CRaDLe shows better generalizability than baseline models.

As can be observed from Tables 5 and 6, one baseline model's extraordinary performance on a specific dataset cannot transfer to other datasets. For example, SelfAttn achieves the best performance among all the baselines on the CodeSearchNet dataset with respect to $R@1$, but perform worse than NeuralBoW on the Code2seq dataset. Comparing with the baselines, CRaDLe presents the best performance on both datasets, which can explicate the good generalizability of CRaDLe.

4.2. Parameter analysis

In this section, we will discuss how the hyperparameters affect the performance of CRaDLe. Three hyperparameters are analyzed, including the number of hidden units in LSTMs, the maximum number of considered statements in the code, and the maximum number of considered tokens in each statement. Figs. 4 and 5 depict the results of the parameter analysis.

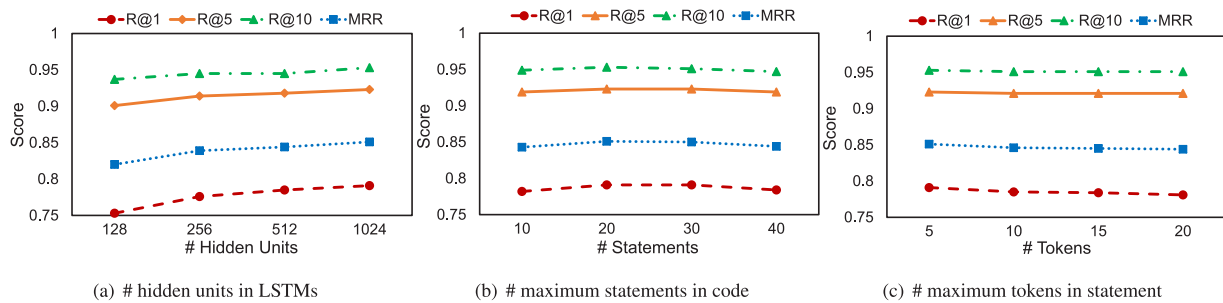


Fig. 4. Parameter sensitivity study for CodeSearchNet.

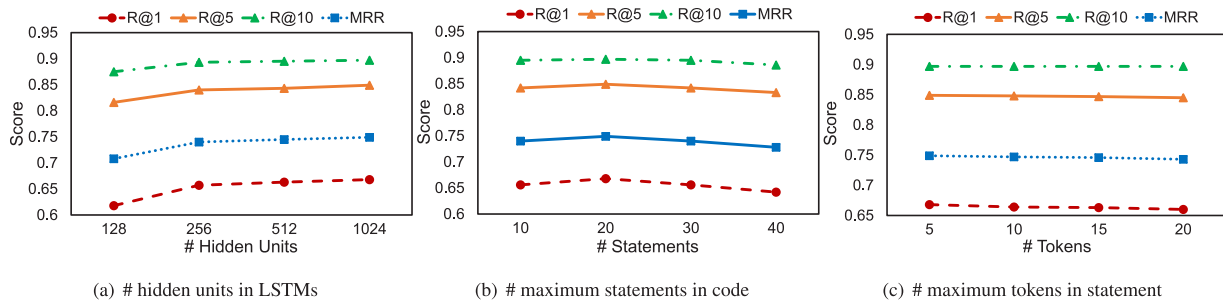


Fig. 5. Parameter sensitivity study for Code2Seq.

4.2.1. # Hidden units in LSTMs

As shown in Figs. 4(a) and 5(a), all the metric values present an increasing trend as the number of hidden units grows. The phenomenon is understandable since more hidden units imply that the model has more parameters to learn and can extract more knowledge from the same input. We can also observe that for each doubling of the number of hidden units, the growth rates of the R@1 scores are 1.9%, 0.54%, 0.41% respectively on the CodeSearchNet dataset. The trend is identical for the Code2Seq dataset. So we can summarize that with an increasing number of the hidden units, the model performance would increase but the increasing rates show a declining tendency. Due to the limitation of the computing source and the marginal enhancement when the number of hidden units is larger than 1024, we choose 1024 as the number of hidden units for our experiment.

4.2.2. # Maximum statements in code

Figs. 4(b) and 5(b) illustrate the variations of the model performance as the maximum number of considered statements increases. We can observe that the metrics achieve the highest values when the number equals 20 and manifests a declining trend as the statement number further increases. As can be found in Tables 3 and 4, the median numbers of the statements in both CodeSearchNet and Code2seq datasets are 7, with the average at around 10. Thus, more statements considered would not be beneficial for capturing the code semantics for most code snippets. In the experiment, we set the maximum number of considered statements in code as 20.

4.2.3. # Maximum tokens in statement

The impact of different maximum numbers of involved tokens in a statement is shown in Figs. 4(c) and 5(c). We can find that when the involved token number increases, the performance presents a downward trend. According to Tables 3 and 4, the average number of tokens in the statements is ~ 3 . So with more tokens recognized, the model could not learn more knowledge of the code snippets. In the experiment, to balance the model performance with the number of tokens considered, we define the maximum number of the tokens in a statement as 5.

```

1 def logs(self, prefix='worker'):
2     logs = []
3     logs += [('success_rate', np.mean(self.success_history))]
4     if self.compute_Q:
5         logs += [('mean_Q', np.mean(self.Q_history))]
6     logs += [('episode', self.n_episodes)]
7
8     if prefix != '' and not prefix.endswith('/'):
9         return [(prefix + '/' + key, val) for key, val in logs]
10    else:
11        return logs

```

Code Listing 2: Successful case 1.

4.3. Ablation study

In the ablation study, we validate the contribution of data dependency or control dependency to CRaDLe and the effectiveness of combining both dependency types. Tables 7 and 8 show the results of the ablation study on the datasets of CodeSearchNet and Code2seq, respectively. CRaDLe_{Full} represents the model utilizes both data dependency and control dependency, CRaDLe_{DataDependency} represents the model only employs data dependency and CRaDLe_{ControlDependency} represents the model only utilizes control dependency.

From the results, we can find that the performance of the model that only utilizes data dependency is very close to the performance of the model with only control dependency, which shows that the importance of data dependency and control dependency is relatively equivalent under our implementation. However, we can find that the model that contains both data dependency and control dependency outperforms the model that only contains one dependency type, especially in terms of the R@1 metric. The results indicate that the combination of data dependency and control dependency is beneficial for effective code search.

4.4. Case studies

Listing 2 shows our predicted code snippet for the query “Generates a dictionary that contains all collected statistics”. We

Table 7
Ablation study on the CodeSearchNet dataset.

Approach	R@1	R@5	R@10	MRR
CRaDLe _{Full}	0.791	0.923	0.951	0.843
CRaDLe _{DataDependency}	0.779	0.910	0.946	0.840
CRaDLe _{ControlDependency}	0.785	0.918	0.950	0.845

Table 8
Ablation study on the Code2seq dataset.

Approach	R@1	R@5	R@10	MRR
CRaDLe _{Full}	0.668	0.849	0.897	0.749
CRaDLe _{DataDependency}	0.645	0.827	0.880	0.724
CRaDLe _{ControlDependency}	0.645	0.828	0.882	0.730

can find that our predicted result correctly matches the given query. Although no overlapping words exist between the code and query, CRaDLe could capture that the code tokens such as `rate` and `compute` are semantically related to the query word “statistics”. Besides, since the semantically-related tokens mainly appear at the line 3, 4 and 5, and do not span the entire code, we guess that the involved dependency information helps to establish the relationships among the statements.

Listing 3 shows another predicted code snippet that accurately matches the given query “Tile N images into one big PxQ image (P,Q)”. Clearly, the function name contains the keywords in the query, e.g., “tile” and “images”. Moreover, the core idea of this query is to tile N images into one image, essentially related to matrix operations. As shown in the Listing 3, the code contains tokens associated with matrix transformation such as `reshape` and `transpose`. So with statement-level tokens explicitly incorporated, CRaDLe could well catch the code functionality.

```

1 def tile_images(img_nhwc):
2     img_nhwc = np.asarray(img_nhwc)
3     N, h, w, c = img_nhwc.shape
4     H = int(np.ceil(np.sqrt(N)))
5     W = int(np.ceil(float(N)/H))
6     img_nhwc = np.array(list(img_nhwc) + [img_nhwc[0]*0 for _ in range(N,
7     ← H*W)])
8     img_Hhwhc = img_nhwc.reshape(H, W, h, w, c)
9     img_HhWhwc = img_Hhwhc.transpose(0, 2, 1, 3, 4)
10    img_Hh_Ww_c = img_HhWhwc.reshape(H*h, W*w, c)
11    return img_Hh_Ww_c

```

Code Listing 3: Successful case 2.

Overall, the above two examples indicate that CRaDLe can accurately capture the code semantics with the statement-level dependency and semantic information integrated.

4.5. Error analysis

Although most of the time, our model returns correct code snippets, we still notice that our model fails under the following two particular circumstances.

4.5.1. Code containing complex mathematical logic

Listing 4 provides a failure case where the code contains complex mathematical logic. The description corresponding to the code is “Convert directly the matrix from Cartesian coordinates (the origin in the middle of image) to Image coordinates (the origin on the top-left of image)”, which includes some mathematical concepts such as “Cartesian coordinates”. Nevertheless, no words related to the mathematical concepts appear in the code. Less knowledge learnt about the mathematical terminology renders the model harder to capture the semantic relevance between the code and

natural language. Future work can incorporate external knowledge such as API documentation or Wikipedia for enhancing the understanding of the mathematical concepts.

```

1 def transform_matrix_offset_center(matrix, y, x):
2     o_x = (x - 1) / 2.0
3     o_y = (y - 1) / 2.0
4     offset_matrix = np.array([[1, 0, o_x], [0, 1, o_y], [0, 0, 1]])
5     reset_matrix = np.array([[1, 0, -o_x], [0, 1, -o_y], [0, 0, 1]])
6     transform_matrix = np.dot(np.dot(offset_matrix, matrix), reset_matrix)
7     return transform_matrix

```

Code Listing 4: Failure case 1.

4.5.2. Code containing function invocation

We also find that the proposed model may fail to capture the code semantics when the code involves function invocation but the details of the invoked function are missing. Listing 5 illustrates such an example, and the corresponding description is “Get successor to key, raises KeyError if a key is max key or key does not exist”. As can be seen in the code example, the execution results strongly rely on the invoked function `succ_item()`, however, the implementation of the invoked function is not detailed. For the case, the code semantics is difficult to be fully captured by the model, leading to failure.

```

1 def succ_key(self, key, default=_sentinel):
2     item = self.succ_item(key, default)
3     return default if item is default else item[0]

```

Code Listing 5: Failure case 2.

5. Discussion

5.1. Dependency embedding approach

In this section, we design another method for representing the dependency information. Specifically, we enrich the dependency matrix with the semantics of the tokens at statement level. The statement-level dependency embedding is calculated as below:

$$\mathbf{p}_i = \frac{\sum_j \mathbf{t}_j v_{ij}}{\max(1, \sum_j v_{ij})}, \forall i = 1, 2, \dots, l, \quad (10)$$

$$\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_l],$$

where \mathbf{t}_j represents the statement-level token embedding for j th statement, which is calculated via Eq. (1). v_{ij} indicates whether the i th statement has a data/control dependency on the j th statement and \mathbf{p}_i is the new dependency embedding.

We evaluated the performance of new dependency embedding methods on the datasets of CodeSearchNet and Code2seq, as shown in Table 9.

From the table, we can find that the new strategy for encoding the dependency information outperforms our original approach in terms of the R@1 and MRR metrics for both datasets. The results indicate that the new approach for the dependency embedding may be more effective than the original approach for the task.

Graph neural networks (GNNs) is also a potential way to represent the dependency between different statements in one code snippet. However, using GNNs for representing the semantic dependency of code is beyond the scope of the work, since the assumption of GNNs that adjacent nodes share similar semantics no long holds for the control dependency information, and it would be more challenging to encode the semantic dependency information through GNNs. In the future, we will investigate various strategies to embed the semantic dependency with GNNs (Li, Ma, Wang, & Zhuang, 2020; Wang et al., 2020).

Table 9

Comparison results with our original models. The best results are highlighted in **bold** fonts.

Dataset	Approach	R@1	R@5	R@10	MRR
CodeSearchNet	CRaDLe _{original}	0.791	0.923	0.951	0.843
	CRaDLe _{new}	0.794	0.920	0.949	0.851
Code2seq	CRaDLe _{original}	0.668	0.849	0.897	0.749
	CRaDLe _{new}	0.676	0.852	0.899	0.756

6. Related work

The work is inspired by the studies related to both code search and code semantics representation learning.

6.1. Code search

In software development, developers accomplish the goal of effective and high quality code by reusing the existing huge amount of available code resources. Prior work has explored a number of methods to find the implicit connections between human language queries and code databases. Early studies concentrate mainly on extracting useful features from both codes and queries. For example, the work (Shepherd, Fry, Hill, Pollock, & Vijay-Shanker, 2007) extracts scattered verbs from queries and applies action-oriented identifier graph model to inspect the result graph, which helps to optimize the queries. In Lu, Sun, Wang, Lo, and Duan (2015), Lu et al. reformulate and extract natural language phrases from source code identifiers since the synonyms in source codes and NL queries may affect the code search result significantly. The work (McMillan et al., 2011) proposes Portfolio uses random surfer to model the navigation behavior of programmers. Then with association model based on Spreading Activation Network (Crestani, 1997), functional relevant functions can be set in the same list. Ponzanelli et al. (Ponzanelli, Bavota, Penta, Oliveto, & Lanza, 2014) propose to retrieve pertinent discussions from Stack Overflow when given a context in the IDE, which saves developers' time spent on formulating more standardized queries.

With rapid development of deep learning, an increasing amount of work has focused on using neural networks for effective code search. In the work (Sachdev et al., 2018), Sachdev et al. first develop neural code search model called NCS to conduct NL search directly over large source code corpora. In Liu et al.' work (Liu, Kim, Murali, Chaudhuri, & Chandra, 2019), they present a neural model called NQE, which expands the queries and improves performance for shorter queries. Codenn embeds both code snippets and natural language descriptions into a unified vector space, in such way that code and its corresponding NL description have similar vectors (Gu et al., 2018). Iyer et al. (Iyer, Konstas, Cheung, & Zettlemoyer, 2016) use attentional long short term memory (LSTM) networks to focus on more cardinal parts of the source code to produce search queries. Cambroner et al. propose UNIF, a bag-of-words-based network which includes API sequences, method tokens, method body tokens and docstring tokens for representing source code (Cambroner et al., 2019). Yao, Peddemail, and Sun regard code annotation and code search as dual task and consider the generated code annotations for better code search (Yao, Peddemail, & Sun, 2019). Husain et al. explore the semantic representations of different neural architectures and they find that self-attention-based architectures achieve the best performance (Husain et al., 2019).

6.2. Representation learning for source code

Prior work has conducted many investigations to effectively represent the semantics of source code. Early studies widely use

machine learning and traditional information retrieval methods. For instance, in Vásquez, McMillan, Poshyvanyk, and Grechanik (2014), Vásquez et al. adopt SVM to discriminate semantic similarities between code snippets and properly categorize software repositories. In the work Kamiya, Kusumoto, and Inoue (2002), programs are morphed into token sequences for facilitating potential code clone detection. Recent work employs deep learning techniques for code semantics learning. Mou et al. adopt tree-structured convolutional neural network (Tree-CNN) to convert source code into distributed vector for program classification (Mou, Li, Zhang, Wang, & Jin, 2016). The work Akbar and Kak (2019) suggests that the order of the embedded words can affect the accuracy of semantic representations. Besides the plain textual information, many studies utilize the structural features of source code, such as abstract syntax tree (AST) and control flow graph (CFG), to enrich the representations of source code. For example, in the work Zhang et al. (2019), AST-based neural network is proposed to capture the structural information of source code. In another work (Chan, Cheng, & Lo, 2012), an API graph and a greedy subgraph search algorithm are utilized to help find the usage of source code, which excavates more semantic details in source code. Functions repetitively called and variables with the same names are involved into Graph Neural Networks (GNNs) for better representing the graph information in the work (Allamanis, Brockschmidt, & Khademi, 2018). Wan et al. propose a multi-modal attention network to combine the heterogeneous sources including AST, CFG and sequences of code tokens (Wan et al., 2019).

7. Conclusions

In this paper, we propose a novel deep neural network named CRaDLe for code retrieval. According to our knowledge, CRaDLe is the first deep learning model which utilizes the program dependency information for the task. CRaDLe learns the code representations with the semantic dependency information combined. Specifically, the dependency information and statement-level tokens are jointly embedded for learning code semantics. Finally, CRaDLe learns unified representations for both code and natural language queries. The experiment results have shown that CRaDLe outperforms the state-of-the-art approaches and the semantic dependency learning is helpful for effective code retrieval.

In the future, we will make a further exploration of the code structure and explicitly incorporate external knowledge such as API documentation to find a better way of representing source code semantics.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This work was supported by the Research Grants Council of the Hong Kong Special Administrative Region, China (CUHK 14210717), the National Natural Science Foundation of China under project No. 62002084, Australian Research Council (ARC) Discovery Project DP200102940, and a key program of fundamental research from Shenzhen Science and Technology Innovation Commission (No. ZX20210035).

References

- Akbar, S. A., & Kak, A. C. (2019). SCOR: source code retrieval with semantics and order. In M. D. Storey, B. Adams, & S. Haiduc (Eds.), *Proceedings of the 16th international conference on mining software repositories* (pp. 1–12). IEEE / ACM, <http://dx.doi.org/10.1109/MSR.2019.00012>, <https://doi.org/10.1109/MSR.2019.00012>.
- Allamanis, M., Brockschmidt, M., & Khademi, M. (2018). Learning to represent programs with graphs. In *6th international conference on learning representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, conference track proceedings*. OpenReview.net.
- Alon, U., Brody, S., Levy, O., & Yahav, E. (2019). Code2seq: Generating sequences from structured representations of code. In *7th international conference on learning representations*. OpenReview.net.
- Brandt, J., Guo, P. J., Lewenstein, J., Dontcheva, M., & Klemmer, S. R. (2009). Two studies of opportunistic programming: interleaving web foraging, learning, and writing code. In D. R. O. Jr., R. B. Arthur, K. Hinckley, M. R. Morris, S. E. Hudson, & S. Greenberg (Eds.), *Proceedings of the 27th international conference on human factors in computing systems* (pp. 1589–1598). ACM.
- Cambroner, J., Li, H., Kim, S., Sen, K., & Chandra, S. (2019). When deep learning met code search. In M. Dumas, D. Pfahl, S. Apel, & A. Russo (Eds.), *Proceedings of the ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering* (pp. 964–974). ACM.
- Chan, W., Cheng, H., & Lo, D. (2012). Searching connected API subgraph via text phrases. In W. Tracz, M. P. Robillard, & T. Bultan (Eds.), *20th ACM SIGSOFT symposium on the foundations of software engineering* (p. 10). ACM.
- Cho, K., van Merriënboer, B., Bahdanau, D., & Bengio, Y. (2014). On the properties of neural machine translation: Encoder-decoder approaches. In D. Wu, M. Carpuat, X. Carreras, & E. M. Vecchi (Eds.), *Proceedings of SSTS@EMNLP 2014, eighth workshop on syntax, semantics and structure in statistical translation, Doha, Qatar, 25 October 2014* (pp. 103–111). Association for Computational Linguistics.
- Crestani, F. (1997). Application of spreading activation techniques in information retrieval. *Artificial Intelligence Review*, 11(6), 453–482.
- Ferrante, J., Ottenstein, K. J., & Warren, J. D. (1987). The program dependence graph and its use in optimization. *ACM Transactions on Programming Languages and Systems*, 9(3), 319–349.
- Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep learning (Vol. 1)*. MIT press Cambridge.
- Gu, X., Zhang, H., & Kim, S. (2018). Deep code search. In M. Chaudron, I. Crnkovic, M. Chechik, & M. Harman (Eds.), *Proceedings of the 40th international conference on software engineering* (pp. 933–944). ACM.
- Henderson, T. A. D., & Podgurski, A. (2016). Sampling code clones from program dependence graphs with GRAPLE. In *Proceedings of the 2nd international workshop on software analytics* (pp. 47–53).
- Husain, H., Wu, H., Gazit, T., Allamanis, M., & Brockschmidt, M. (2019). Code-searchnet challenge: Evaluating the state of semantic code search. CoRR [arXiv:abs/1909.09436](https://arxiv.org/abs/1909.09436).
- Iyer, S., Konstas, I., Cheung, A., & Zettlemoyer, L. (2016). Summarizing source code using a neural attention model. In *Proceedings of the 54th annual meeting of the association for computational linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long papers*. The Association for Computer Linguistics.
- Kamiya, T., Kusumoto, S., & Inoue, K. (2002). Cfinder: A multilingual token-based code clone detection system for large scale source code. *IEEE Transactions on Software Engineering*, 28(7), 654–670.
- Kim, Y. (2014). Convolutional neural networks for sentence classification. In A. Moschitti, B. Pang, & W. Daelemans (Eds.), *Proceedings of the 2014 conference on empirical methods in natural language processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, a meeting of SIGDAT, a special interest group of the ACL* (pp. 1746–1751). ACL.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. Y. Bengio, Y. LeCun (Eds.), *3rd international conference on learning representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, conference track proceedings*.
- Lee, C.-Y., Gallagher, P. W., & Tu, Z. (2016). Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree. In *Artificial intelligence and statistics* (pp. 464–472).
- Li, M., Ma, Z., Wang, Y. G., & Zhuang, X. (2020). Fast haar transforms for graph neural networks. *Neural Networks*, 128, 188–198.
- Li, Y., Wang, S., Nguyen, T. N., & Nguyen, S. V. (2019). Improving bug detection via context-based code representation learning and attention-based neural networks. *Proceedings of ACM Programming Languages*, (OOPSLA), 162:1–162:30.
- Lin, Z., Feng, M., dos Santos, C. N., Yu, M., Xiang, B., Zhou, B., et al. (2017). A structured self-attentive sentence embedding. In *5th international conference on learning representations, ICLR 2017, Toulon, France, April 24-26, 2017, conference track proceedings*. OpenReview.net.
- Liu, J., Kim, S., Murali, V., Chaudhuri, S., & Chandra, S. (2019). Neural query expansion for code search. In T. Mattson, A. Muzahid, & A. Solar-Lezama (Eds.), *Proceedings of the 3rd ACM SIGPLAN international workshop on machine learning and programming languages* (pp. 29–37). ACM.
- Lu, M., Sun, X., Wang, S., Lo, D., & Duan, Y. (2015). Query expansion via wordnet for effective code search. In Y. Guéhéneuc, B. Adams, & A. Serebrenik (Eds.), *22nd IEEE international conference on software analysis, evolution, and reengineering* (pp. 545–549). IEEE Computer Society.
- Lv, F., Zhang, H., Lou, J., Wang, S., Zhang, D., & Zhao, J. (2015). Codehow: Effective code search based on API understanding and extended boolean model (e). In M. B. Cohen, L. Grunske, & M. Whalen (Eds.), *30th IEEE/ACM international conference on automated software engineering* (pp. 260–270). IEEE Computer Society.
- McMillan, C., Grechanik, M., Poshyvanyk, D., Xie, Q., & Fu, C. (2011). Portfolio: finding relevant functions and their usage. In R. N. Taylor, H. C. Gall, & N. Medvidovic (Eds.), *Proceedings of the 33rd international conference on software engineering* (pp. 111–120). ACM.
- Mou, L., Li, G., Zhang, L., Wang, T., & Jin, Z. (2016). Convolutional neural networks over tree structures for programming language processing. In D. Schuurmans, & M. P. Wellman (Eds.), *Proceedings of the thirtieth AAAI conference on artificial intelligence* (pp. 1287–1293). AAAI Press.
- Ponzanelli, L., Bavota, G., Penta, M. D., Oliveto, R., & Lanza, M. (2014). Mining stackoverflow to turn the IDE into a self-confident programming prompter. In P. T. Devanbu, S. Kim, & M. Pinzger (Eds.), *11th working conference on mining software repositories, MSR 2014, Proceedings, May 31 - June 1, 2014, Hyderabad, India* (pp. 102–111). ACM.
- PYPL (2020). PYPL popularity of programming language. <http://pypl.github.io/PYPL.html>.
- Sachdev, S., Li, H., Luan, S., Kim, S., Sen, K., & Chandra, S. (2018). Retrieval on source code: a neural code search. In J. Gottschlich, & A. Cheung (Eds.), *Proceedings of the 2nd ACM SIGPLAN international workshop on machine learning and programming languages*, (pp. 31–41). ACM.
- Shepherd, D. C., Fry, Z. P., Hill, E., Pollock, L. L., & Vijay-Shanker, K. (2007). Using natural language program analysis to locate and understand action-oriented concerns. In B. M. Barry, & O. de Moor (Eds.), *ACM International Conference Proceeding Series: 208, Proceedings of the 6th international conference on aspect-oriented software development* (pp. 212–224). ACM.
- Vásquez, M. L., McMillan, C., Poshyvanyk, D., & Grechanik, M. (2014). On using machine learning to automatically classify software applications into domain categories. *Empirical Software Engineering*, 19(3), 582–618.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). Attention is all you need. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, R. Garnett (Eds.), *Advances in neural information processing systems 30: Annual conference on neural information processing systems 2017* (pp. 5998–6008).
- Wan, Y., Shu, J., Sui, Y., Xu, G., Zhao, Z., Wu, J., et al. (2019). Multi-modal attention network learning for semantic source code retrieval. In *34th IEEE/ACM international conference on automated software engineering* (pp. 13–25). IEEE.
- Wang, Y., Li, M., Ma, Z., Montúfar, G., Zhuang, X., & Fan, Y. (2020). Haar graph pooling. In *Proceedings of Machine Learning Research: 119, Proceedings of the 37th international conference on machine learning, ICML 2020, 13-18 July 2020, virtual event* (pp. 9952–9962). PMLR.
- Wang, S. L., & Manning, C. D. (2012). Baselines and bigrams: Simple, good sentiment and topic classification. In *The 50th annual meeting of the association for computational linguistics, proceedings of the conference, July 8-14, 2012, Jeju Island, Korea - Volume 2: Short papers* (pp. 90–94). The Association for Computer Linguistics.
- Yao, Z., Peddamail, J. R., & Sun, H. (2019). Coacor: Code annotation for code retrieval with reinforcement learning. In L. Liu, R. W. White, A. Mantrach, F. Silvestri, J. J. McAuley, R. Baeza-Yates, & L. Zia (Eds.), *The world wide web conference* (pp. 2203–2214). ACM.
- Yin, P., & Neubig, G. (2017). A syntactic neural model for general-purpose code generation. In R. Barzilay, & M. Kan (Eds.), *Proceedings of the 55th annual meeting of the association for computational linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long papers* (pp. 440–450). Association for Computational Linguistics.
- Zhang, J., Wang, X., Zhang, H., Sun, H., Wang, K., & Liu, X. (2019). A novel neural source code representation based on abstract syntax tree. In J. M. Atlee, T. Bultan, & J. Whittle (Eds.), *Proceedings of the 41st international conference on software engineering* (pp. 783–794). IEEE / ACM.