
DEGRAPHCS: EMBEDDING VARIABLE-BASED FLOW GRAPH FOR NEURAL CODE SEARCH

Chen Zeng

School of Computer

National University of Defense Technology

zengchen15@nudt.edu.cn

Yue Yu

School of Computer

National University of Defense Technology

yuyue@nudt.edu.cn

Shanshan Li

School of Computer

National University of Defense Technology

shanshanli@nudt.edu.cn

Xin Xia

School of Computer

Monash University

xin.xia@monash.edu

Zhiming Wang

School of Computer

National University of Defense Technology

wangzhiming14@nudt.edu.cn

Mingyang Geng

School of Computer

National University of Defense Technology

gengmingyang13@nudt.edu.cn

Linxiao Bai

School of Computer

National University of Defense Technology

linxiao_b@nudt.edu.cn

Wei Dong

School of Computer

National University of Defense Technology

wdong@nudt.edu.cn

Xiangke Liao

School of Computer

National University of Defense Technology

xkliao@nudt.edu.cn

March 25, 2021

ABSTRACT

With the rapid increase in the amount of public code repositories, developers maintain a great desire to retrieve precise code snippets by using natural language. Despite existing deep learning based approaches (*e.g.*, DeepCS and MMAN) have provided the end-to-end solutions (*i.e.*, accepts natural language as queries and shows related code fragments retrieved directly from code corpus), the accuracy of code search in the large-scale repositories is still limited by the code representation (*e.g.*, AST) and modeling (*e.g.*, directly fusing the features in the attention stage). In this paper, we propose a novel learnable *deep Graph for Code Search* (called **DEGRAPHCS**), to transfer source code into variable-based flow graphs based on the intermediate representation technique, which can model code semantics more precisely compared to process the code as text directly or use the syntactic tree representation. Furthermore, we propose a well-designed graph optimization mechanism to refine the code representation, and apply an improved gated graph neural network to model variable-based flow graphs. To evaluate the effectiveness of **DEGRAPHCS**, we collect a large-scale dataset from GitHub containing 41,152 code snippets written in C language, and reproduce several typical deep code search methods for comparison. Besides, we design a qualitative user study to verify the practical value of our approach. The experimental results have shown that **DEGRAPHCS** can

achieve state-of-the-art performances, and accurately retrieve code snippets satisfying the needs of the users.

Keywords Code search, graph neural networks, intermediate representation, deep learning

1 Introduction

Code search has received increasing attention in recent years [1, 2, 3, 4, 5]. The goal of code search is to retrieve code fragments which best meet developers' need by performing natural language queries over a large code corpus. With the availability of immense and rapidly growing source code repositories such as GitHub and Stack Overflow, it is more convenient for developers to search the needed code with certain functionality and reuse it in their own programs. However, increasingly complex and diverse code implementations also bring great challenges to perform a precise code search.

In the early stage, code search approaches were proposed on the basis of information retrieval techniques, especially key words matching mechanism [6, 7, 8, 9, 10, 11, 12, 13, 14]. However, a common problem of these works is the lack of structural or semantic information from the source code since they simply consider code and queries as plain texts. Recently, deep learning technologies have been applied to represent code and queries for code search [1, 2, 3, 4, 5] to tackle the above issues. The typical approach, called DeepCS [2], presented a code search engine by learning a joint-embedding of a method description and its corresponding code snippet. Moreover, Wan et al. [5] designed a Multi-Modal Attention Network (MMAN) to capture various code features simultaneously, such as code tokens, abstract syntactic tree (AST) and statement-based control-flow graph (S-CFG).

However, existing deep learning based approaches are still limited from two major aspects. First, in reality, code with different syntax may achieve the same functionality, while code with similar structural features may express totally different code semantics. Thus, the token (*e.g.*, method name or identifiers) and structural features (*e.g.*, AST or S-CFG) are hard to precisely express the in-depth semantics of source code in various forms (as shown in Fig. 1 and Fig. 2). Second, existing methods cannot fully exploit multiple valuable features extracted from the source code. In concrete, some models did not fuse different modalities of source code effectively, which does not bring much improvement yet increases the complexity instead. For example, MMAN proposed an attention network to assign learnable weights to three different modality of source code (*i.e.*, Token + AST + S-CFG), while only outperforming the single token-based modality by 4.63% and 6.12% in terms of MRR (mean reciprocal rank) and SuccessRate@1.

These aforementioned limitations inspire us to design a model to effectively integrate several deep semantic information and learn a precise code representation. In our work, we explore a novel code representation based on data and control flow extracted from LLVM IR (Intermediate Representation) [15], one type of intermediate code acquired from source code. Different from the existing statement-based data and control flow representation method [16], we refine the construction of the variable-based flow graph to better describe the dependencies between code variables. In concrete, the nodes in the graph represent the tokens appeared in LLVM IR and the edges represent the data and control dependencies between the tokens. Furthermore, we design an optimization mechanism while modeling the graph to remove the redundant information brought by LLVM IR without changing the semantics. Finally, we employ an attentional gated graph neural network to embed the flow graph into a high-dimensional vector space to further perform code search tasks. Through this procedure, multiple semantic features of code, *i.e.*, tokens, variable-based data and control flow, can be simultaneously represented and accurately express the deep semantics of code.

To evaluate the effectiveness of our proposed model, we collect our dataset from GitHub containing 41,152 code snippets written in C language and perform code search experiments. Experimental results show that DEGRAPHCS improves the top-1 hit rate of code search from 34.05% to 43.05% when compared with the state-of-the-art methods. To simulate the actual code search scenario, we design an online code search tool, which takes 50 practical descriptions randomly chosen from the test set as candidate queries. For each query, 5 experienced participants manually label the relevant results they need returned by our proposed model DEGRAPHCS and three competitive approaches (*i.e.*, DeepCS, UNIF and MMAN). The results of automatic evaluation and manual evaluation both confirm the effectiveness of DEGRAPHCS.

The main contributions of this paper are summarized as follows:

- We propose a novel semantic code representation method called DEGRAPHCS, which can integrate token, data flow and control flow into a variable-based graph more precisely than traditional approaches (*e.g.*, AST). All of our code and data are released at <https://github.com/degraphcs/DeGraphCS>.

- We design a graph optimization mechanism to streamline the representation of graph by reducing 51.88% redundant nodes, which can significantly improve the performance of DEGRAPHCS by 13.77% and 18.92% in terms of MRR and SuccessRate@1.
- We collect a large-scale dataset from GitHub containing 41,152 code snippets written in C language, and reproduce several competing code search models to make comparison.
- We conduct experiments on the trained models, and the results of automatic evaluation demonstrate that DEGRAPHCS outperforms the state-of-the-art method (*i.e.*, MMAN) by 19.14% and 26.43% in terms of MRR and SuccessRate@1. Besides, DEGRAPHCS achieves the best performance in our qualitative user study.

The remainder of this paper is organized as follows. In Section 2, we present two motivating examples. In Section 3, we first provide an overview of our proposed model and then describe the details of each part in our model. In Section 4, we elaborate the experimental setup and report the experimental results. In Section 5, we briefly review the related works. Finally, in Section 6, we conclude our study and future work.

2 Motivating Example

In this section, we show two motivating examples in Fig. 1 and Fig. 2 to illustrate the advantages of our semantic feature-based code representation approach over the existing methods which are based on structural features. Here, we argue that code snippets with same functionality may have different implementation, while code snippets with totally different code semantics may have similar structural features. Therefore, we need to precisely represent the source code on the semantics for performing a satisfying code search task. In other words, we need to find a more precise representation for source code, making semantically similar code have similar representation.

Fig. 1(d-f) show three simple C code snippets as example, all of them aim to sum the values in the specified array. Fig. 1(a-c) represent the corresponding ASTs of the three code snippets. Our variable-based flow graph constructed from LLVM IR is illustrated in Fig. 1(g) and Fig. 1(h). Here, nodes in the flow graph are tokens of the intermediate code (*e.g.*, generated by LLVM IR) and the edges either represent data dependencies (shown in solid line of light blue) or control dependencies (shown in dotted line of red). From Fig. 1(d-f), we can clearly see that the three code examples implement the same function but have different writing format (“for”, “while” and recursive loop separately), which result in totally different syntactic structure in Fig. 1(a-c). The difference between the three ASTs is highlighted in red and yellow. In concrete, in Fig. 1(a) and Fig. 1(c), the parent node of the sub-tree connected in red line is “get_sum”, while in Fig. 1(b), the corresponding parent node is “for” instead. Moreover, in Fig. 1(c), “if” node’s sub-tree connected in yellow line is completely different from the sub-tree of “while” node and “for” node in Fig. 1(a) and Fig. 1(b). However, as shown in Fig. 1(g-h), when exploiting our variable-based flow graph, the three code snippets with the same semantics are represented almost the same. In concrete, except for node “ptr_1” and “ptr”, every node in Fig. 1(g) can find corresponding node in Fig. 1(h). And compared with node in Fig. 1(g), the corresponding node have same control dependency and data dependency (*e.g.*, in Fig. 1(g) and Fig. 1(h), two “return” nodes both depend on “sum_1” and “label_false”). In fact, in Fig. 1(g), both “ptr_1” and “ptr” refer to variable “ptr”, their dependencies are exactly the same as the “ptr” in Fig. 1(h). In short, almost every data and control dependency in Fig. 1(g) can find corresponding one in Fig. 1(h). And the data and control dependency imply code semantics. Thus Fig. 1(g) and Fig. 1(h) can be regarded as same on code semantics. In all of the above, our flow graph establish an unifying representation for semantically similar code. And the unifying representation relies on the fact that data and control dependencies usually capture deeper code semantics than shallow syntactic features. Therefore, we argue that code snippets with similar semantics should maintain similar code representations.

Furthermore, to better explain our motivation that code snippets expressing different semantics should differ in code representations, we construct two segments of code snippets (shown in Fig. 2(c) and Fig. 2(d)) as the second example. Fig. 2(a)(b) and Fig. 2(e)(f) depict the AST and our variable-based flow graph of each code snippet. As shown in Fig. 2(c) and Fig. 2(d), although the two code examples consist of exactly the same statements, the different order will lead to different functionalities while running. From Fig. 2(a) and Fig. 2(b), we can see that the two tree structures are the same except for the relative positions of the two sub-trees of the root node (connected in red/yellow respectively), and the existing works like MMAN will finally establish the same representation when applying Tree-LSTM [17]. However, the corresponding variable-based flow graphs differ significantly in Fig. 2(e) and Fig. 2(f) because of the underlying data dependencies between variables. The two examples have shown that our proposed variable-based flow graph representation can represent precise semantics of code while the syntactic structures fail.

Inspired by the above two examples, we conclude that data and control dependencies can complement the drawbacks of the structural features-based code representation methods. However, the existing works extract data and control flows on the basis of statements, obtaining vector representations of code snippets by applying word embedding technique

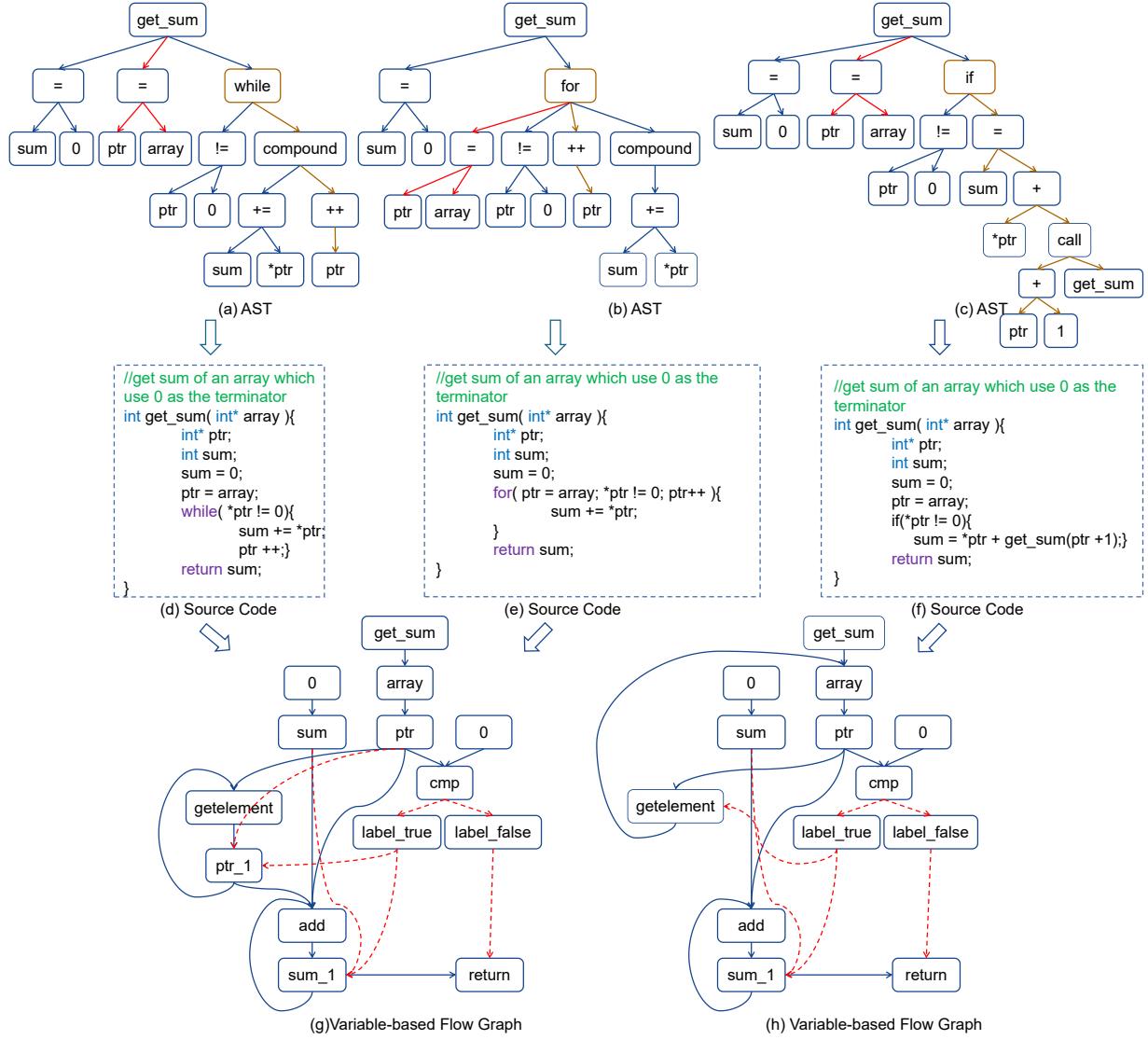


Figure 1: The first illustrative example shows the code snippets with the same semantics and their corresponding ASTs and variable-based flow graphs.

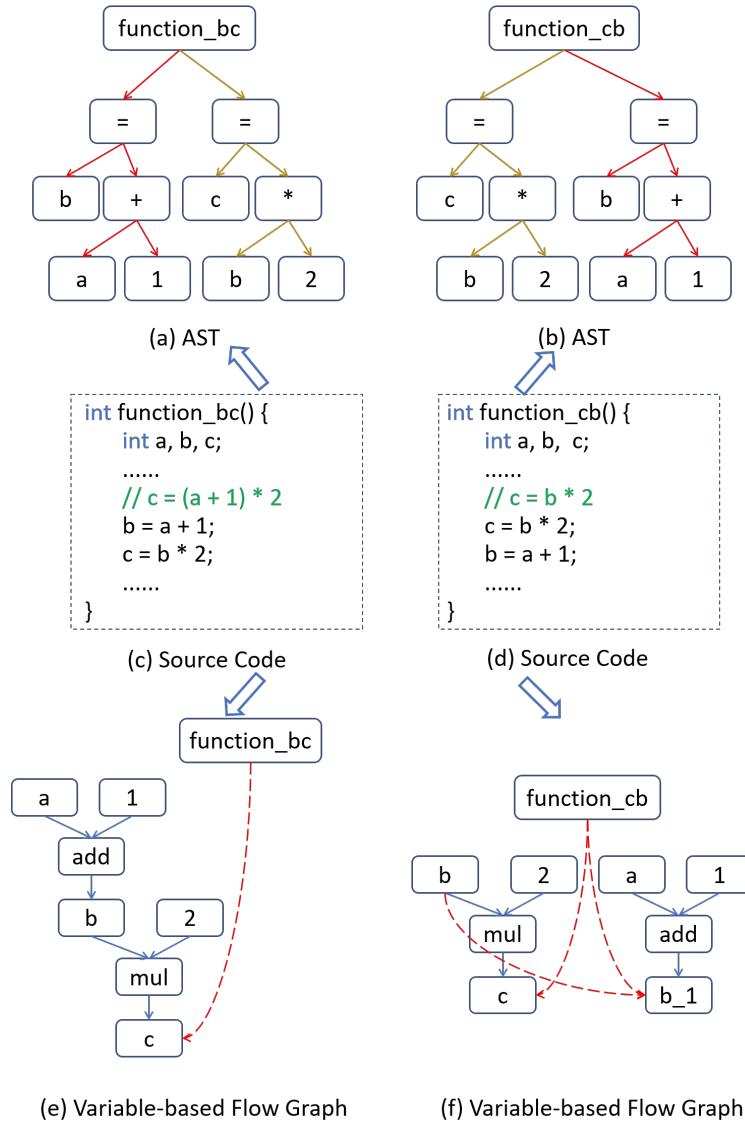


Figure 2: The second illustrative example shows the code snippets with different semantics and their corresponding ASTs and flow-based graphs.

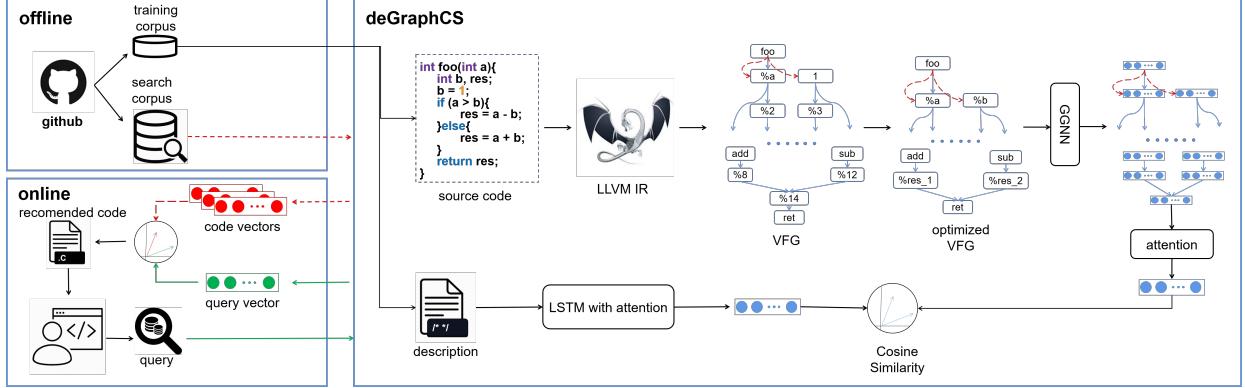


Figure 3: The overall workflow of deGRAPHCS, containing the offline data preparing, the online inference and the network architecture part.

such as skip-gram [16]. A prominent problem is that coarse-grained statements usually cannot accurately capture the correlation between the tokens of code and query. Therefore, we propose our fine-grained variable-based flow graph method to precisely model the relationship between tokens in code snippets.

3 The Proposed Model

In this section, we first introduce an overview of our proposed network architecture. Then, we present our neural code representation mechanism, including background of compilation and LLVM IR, the variable-based flow graph building mechanism and optimization mechanism. Finally, we present our comment description representation and model learning mechanism in detail.

3.1 An Overview

Fig. 3 is an overview of workflow of deGRAPHCS model, which is composed of three parts: the upper left part denotes the offline data preparing process; the bottom left part denotes the online inference process; the right part denotes the details of the network architecture. For the network architecture, deGRAPHCS first embeds the neural code and comment to the vector representations, and then learns the relationship by minimizing the ranking loss function in the training process. We will describe each part of the architecture in the following sections.

3.2 Neural Code Representation

For code representation, we first integrate data dependencies and control dependencies into graphs by analyzing different kinds of LLVM IR instructions. In concrete, we construct the data dependencies based on the address operation instructions (*e.g.*, “load”, “store”) and the computation related constructions (*e.g.*, “add” and “sub”). Besides, the control dependencies are constructed based on the jump instructions (*e.g.*, “br”) and address operation instructions. The goal of code search is to better match the semantics in code with the keywords in queries, excessive information may hinder the model from learning the fine-grained relationship between the source code and queries. Therefore, we propose several mechanisms to optimize the graph with the aim of decreasing the noises in the model training process and improving the training efficiency. Finally, we feed the graph into a GGNN[18] with attention mechanism to learn the vector representation of the code. In order to better explain our proposed neural code representation method, an illustrative example code associated with its IR, and our variable-based flow graph building and optimizing result is shown in Fig. 4.

3.2.1 Compilation and LLVM IR

Most popular compilers, such as LLVM and GCC, support multiple programming languages and hardware targets. With the aim of avoiding duplications in code optimization techniques, the compilers require a strict separation among the source language, IR, and the target machine code which will be mapped to a specific hardware. LLVM IR supports various architectures and can represent optimized code inherently. The IR in LLVM is given in Static Single Assignment (SSA) form [19], which guarantees that every variable is assigned only once. As shown in Fig. 4(b), LLVM divides the IR statements into several blocks represented by the corresponding labels shown in different colors. For the

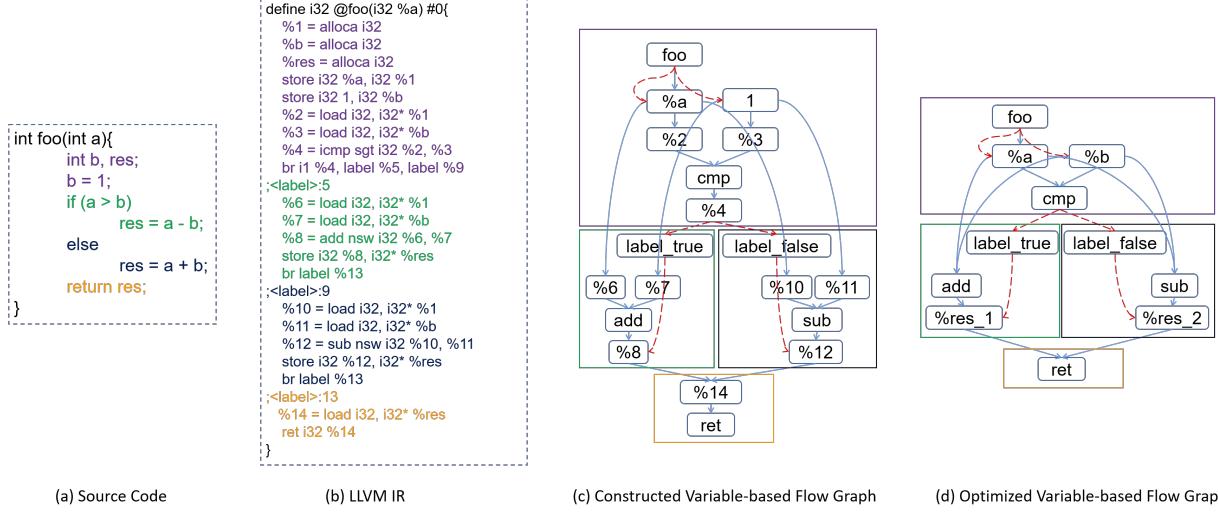


Figure 4: An illustrative example shows a code snippet with its LLVM IR equivalent and our graph building as well as optimizing manner.

instructions, regarding the third line in “label 9”, an instruction (IR statement) in our algorithm is mainly composed of three parts: opcode (“sub”), operand (%10, %11) and result (%12).

Algorithm 1: Variable-based Flow Graph Building Process

Input: LLVM IR (shown in Fig. 4(b)).

Output: the constructed variable-based flow graph (shown in Fig. 4(c)).

- 1: **for** read the instructions of IR by line **do**
 - 2: **Case** Computation instructions:
 - 3: **Case** “call/invoke”:
 - 4: Build an edge parameters → function name.
 - 5: **Default:**
 - 6: Build an edge operands → opcode.
 - 7: Build an edge opcode → result.
 - 8: **Case** Address instructions:
 - 9: **for** each load instruction operated on *addr* (*i.e.*, *a* = load *addr*) **do**
 - 10: *value_list* = SearchCFG(*addr*, *inst*).
 - 11: Build edges *value_list* → *a*.
 - 12: **end for**
 - 13: **for** all instructions “store *x* *addr*” **do**
 - 14: Connect *x* sequentially according to the order in CFG.
 - 15: **end for**
 - 16: **Case** “br” (*i.e.*, br %val, label 1,label2):
 - 17: Build edges condition → labels.
 - 18: Build edges labels → “store” variables.
 - 19: **end for**
-

3.2.2 Building Variable-based Flow Graph (VFG)

To derive a precise semantic representation of source code, we construct the graph at the granularity of variables to capture token information of source code. In concrete, we build data dependencies and control dependencies between variables to capture the data and control flow information from source code. An illustrative example of the variable-based flow graph is shown in Fig. 4(c), which is the initial graph constructed from the LLVM IR shown in Fig. 4(b). The nodes in our graph can be variables, opcode or label identifiers, appearing in the figure as rectangles. Correspondingly, an edge either represents data dependency (in blue solid line), or control dependency (in red dotted line). Given the LLVM IR, the whole graph building process is recorded in **Algorithm 1**. In concrete, we first extract

the identifiers in each IR instruction as nodes. Then, we build data dependencies and control dependencies between nodes according to different types of instructions as follows.

Data dependency. There exist data dependencies in the computation related constructions (*e.g.*, “add”, “sub”) and the address operation instructions (*e.g.*, “load”, “store”). First, we build data dependencies according to instructions which are related to computation such as “add”, “sub” etc. For example, regarding the third line in the second block shown in green of Fig. 4(b) (“%8 = add nsw i32 %6, %7”), we build data dependencies by linking the operands (%6 and %7) to the opcode (“add”), and then linking the opcode (“add”) to the result (%8) shown in the green square of Fig. 4(c). We specially deal with the “call/invoke” instruction since the operands are the parameters of corresponding function. In concrete, when we build graph for current function (*e.g.*, “get_sum”), we treat the function call instruction in two situation. First, if the called function is the external function which is not the current function (“get_sum”), we treat the name as opcode instead of “call/invoke”, and link the operands to the called function name. Second, if the called function is the current function, it is regarded as a recursive call. Thus we regard the node which link to “return” node as the result of called function, and we link the result to the nodes which use the called function. Moreover, we regard the input of called function is passed into the parameter of current function, thus we link the input node to the parameter node. For example, when we construct data dependency for the function “get_sum(array)” in Fig. 1(f), the function “get_sum(array)” call itself recursively (“sum = *ptr + get_sum(ptr +1);”). Thus we deal with the function call “get_sum(ptr +1)” according to the second situation. As Fig. 1(h), we can find “sum_1” is linked to “return” which implies “get_sum” return “sum_1” as result. Thus we link “sum_1” to “add” which use the result of “get_sum” to compute. Then we link “getelement” which denotes the input “ptr +1” to parameter node “array” of function “get_sum(array)”.

Algorithm 2: Search All the Variables from *addr* in *inst*

Input: *addr, inst*.

Output: *value_list*.

```

1: pre_list  $\leftarrow$  the last instructions pointed to inst.
2: for each instruction pre_inst in the pre_list do
3:   if pre_inst has been searched then
4:     continue
5:   else if pre_inst is “store x addr” then
6:     Append the variable x to the value_list.
7:   else
8:     list = SearchCFG(addr, inst)
9:     Append all the values in list to value_list
10:  end if
11: end for
```

Second, we build the data dependencies of variables in “load” and “store” instructions. In concrete, when the variables need to be used, LLVM will load the corresponding values from the allocated address by “load” instruction. Similarly, when variables need to be assigned new values, IR will store the new values into address by “store” instruction. The reason to treat “load” and “store” instructions separately is that in these two instructions, the address may store multiple values from different variables, it is hard to build one-to-one mapping relationship between each address and variable. Therefore, we build data dependencies only between variables. For example, regarding the “load” instructions of the second block shown in green of Fig. 4(b), although the variables %6 and %7 are loaded from the address %1 and %b, the edges are connected from variables %a and 1, which are the true sources. To achieve this, we need to traverse all the “load” instructions, and handle each instruction following the function in **Algorithm 2**.

Control dependency. There exist control dependencies in the jump instructions (*e.g.*, “br”) and the address operation instructions (*e.g.*, “load”, “store”). We exploit address operation instructions since multiple variables stored in an address usually maintain a sequential order. Thus, we complete control dependencies of variables through two aspects as follows:

First, we build control dependencies between variables and the condition identifiers. These conditions appear in condition jump instructions such as “br” instruction as shown in the last purple line of Fig. 4(b). Label identifiers are the entries of basic blocks, and the condition determines which label to jump. Thus, we construct control dependencies by linking the condition to all label identifiers. After that, to make the whole graph connected, we should build control dependencies between the label and the corresponding basic block’s variables. In our algorithm, we link the label identifier to the variables in “store” instruction as shown in the red line from label_true to %8 of Fig. 4(c). The reason is that if a variable updates its value, a new value will be stored in the corresponding address by a store instruction.

Second, multiple assignments of the same variable will generate different variables to be stored in the address and the variables usually maintain a sequential order. Therefore, we build control dependencies between these variables of the same address. As shown in line 13 of Algorithm 1, we need to traverse every store instruction which stores a variable to address \mathbf{x} , and find the next “store” instruction which stores a new variable to address \mathbf{x} , then build a connection from the previous variable to the latter variable until all the variables are connected in sequential order.

3.2.3 Optimizing Variable-based Flow Graph

After constructing the variable-based flow graph, we need to optimize the graph to decrease the noises and improve the training efficiency. The optimizing method is composed of the following four steps.

1. First, in LLVM IR, variables in “store” instruction are named with numbers (e.g., %1, %2). Since the goal of code search is to better match the tokens (e.g., identifiers, functions names) in code with the keywords in queries, excessive numbers may hinder the model from filtering the critical information of the flow graph. Therefore, we replace the numbers with the corresponding variable names.
2. Second, plenty of opcode are too trivial to represent the code semantics needed for code search. Therefore, we remove these opcode nodes from the graph by linking the predecessors of the opcode to their successors. In concrete, the trivial opcode can mainly be divided into three kinds. The first kind of opcode are those related to memory access and addressing since they operate on the variable address. The second are opcode related to conversion since they aim to transform the type of the data (e.g., change “int” type to “string” type). The third are opcode related to operations on exceptions.
3. Third, in LLVM IR, many temporary registers are generated to store the intermediate values and these registers have no corresponding variables. Thus, we remove these variable nodes by linking their predecessors to their successors.
4. Fourth, we compress the control flow graph by merging the “isolated” blocks. In our constructed variable-based flow graph, suppose block a is the predecessor of block b , if block a has only one successor and block b has only one predecessor, then we merge the two blocks to compress the control flow.

3.2.4 Graph2Vec

Since our constructed graph is a directed graph with multiple types of edges, we utilize an improved gated graph neural network (GGNN) with attention mechanism to learn the vector representation of the code. GGNN is a neural network architecture for embedding graphs with multiple types of edges. In our graph $G = (V, E)$, V denotes a set of nodes(v, l_v), and E denotes a set of edges($(v_i, v_j, l_{(v_i, v_j)})$). l_v denotes the label of node v which is consisted of the variables in IR instructions. $l_{(v_i, v_j)}$ denotes the label of the edge from v_i to v_j which includes two types: data dependency and control dependency.

GGNN learns the vector representation of G by message passing mechanism as follows. First, we initialize each node $v \in V$ with a one-hot embedding vector(h_v^0) according to l_v . Then, we train the embeddings of all nodes through multiple iterations. In iteration t , each node v_i obtains message $m_t^{v_j \rightarrow v_i}$ from neighbour v_j as: $m_t^{v_j \rightarrow v_i} = W_{l_{(v_i, v_j)}} h_{v_j}^{t-1}$, $W_{l_{(v_i, v_j)}}$ is the weight matrix specified by the type of edges, it map message from neighbour v_j into a shared space. The weight matrix is learned during training process. Then all message from the neighbours of v_i are aggregated in the following equation:

$$m_t^i = \sum_{v_j \in \text{Neighbour}(v_i)} (m_t^{v_j \rightarrow v_i}) \quad (1)$$

Then, GGNN uses GRU (Gated Recurrent Unit)[20] to update the embedding of each node v_i . GRU uses aggregated message and past state $h_{v_i}^{t-1}$ to update current state as $h_{v_i}^t = \text{GRU}(m_t^i, h_{v_i}^{t-1})$. Finally, since different nodes contribute differently to the code semantics, we exploit the attention mechanism to calculate the importance of different nodes. We first allocate weights for each node v_i as:

$$\alpha_i = \text{sigmoid}(f(h_{v_i}) \cdot u_{vfg}) \quad (2)$$

α_i denotes the weight of node v_i , $f(\cdot)$ denotes the linear layer, \cdot denotes the inner project function and u_{vfg} denotes the context vector which is a high level representation of the whole nodes in graph. u_{vfg} is realized as a linear layer which randomly initialized and jointly learned during training. Then, we obtain the embedding of the whole graph h_{vfg} as:

$$h_{vfg} = \sum_{v_i \in V} (\alpha_i h_{v_i}). \quad (3)$$

3.3 Comment Description Representation

For comment representation, we apply LSTM[21] to learn the corresponding representations. The embedding h_i^{des} of each word in comment is calculated as $h_i^{des} = LSTM(h_{i-1}^{des}, w(d_i))$, where $i = 1, \dots, |d|$, $|d|$ denotes the length of the comment description, and w denotes the word embedding layer to embed each word into a vector. Since different parts of the comment make different contributions to the final vector representation, we adopt an attention mechanism [22] to capture the fine-grained relevance between the hidden states and the final comment representation. In concrete, we apply an attention layer to calculate the attention score $\alpha^{des}(i)$:

$$\alpha^{des}(i) = \frac{\exp(f(h_i^{des}) \cdot u^{des})}{\sum_{k=1}^n \exp(f(h_k^{des}) \cdot u^{des})} \quad (4)$$

where \cdot denotes the inner product of h_i^{des} and u^{des} , $f(\cdot)$ denotes a linear layer and u^{des} denotes the context vector which is a high level representation of the whole tokens in comment. The context vector u^{des} is randomly initialized and jointly learned during training. Then, the final representation of the comment description $E_{|d|}^{des}$ can be calculated as:

$$E_{|d|}^{des} = \sum_{i=1}^{|d|} \alpha^{des}(i) h_i^{des} \quad (5)$$

3.4 Model Training

Now we obtain all code representation (C) and description representation (D). To search the code precisely for each query, the model should make the code representation similar with the correct description representation, and make the code representation different with the incorrect description representation. In other words, for a code snippet representation $c \in C$, the associated correct description $d^+ \in D$ and a randomly chosen incorrect description $d^- \in D$, we need to make the vector representation of the pair $< c, d^+ >$ similar and the vector representation of the pair $< c, d^- >$ different. Therefore, we train the model by minimizing the loss function $L(\theta)$ in formulation of:

$$L(\theta) = \sum_{c \in C} \sum_{d^+, d^- \in D} \max(0, \beta - \cos(c, d^+) + \cos(c, d^-)) \quad (6)$$

where θ denotes the model parameters, d^+ denotes the correct description representation, d^- denotes the incorrect description representation. The $\cos(\cdot, \cdot)$ function is to measure the cosine similarity between two vector representations, and β denotes the constant margin.

4 Experiments and Results

In this section, to evaluate the performance of DEGRAPHCS in code search, we perform several experiments to answer the following research questions:

- RQ1: How effective is our proposed DEGRAPHCS?
- RQ2: What is the effectiveness of our natural integration of multiple information, *i.e.*, code tokens, variable-based data and control flow extracted from IR in DEGRAPHCS when compared with the existing attention-based multi-modal representation method?
- RQ3: How does our graph optimizing mechanism affect the final retrieval performance?
- RQ4: What is the performance of DEGRAPHCS when varying the comment length, code length, VFG nodes number and the longest path length of VFG?
- RQ5: What is the performance of DEGRAPHCS for helping developers in actual programming?

RQ1 is to investigate whether DEGRAPHCS outperforms the state-of-the-art deep code search models. RQ2 is to evaluate the effectiveness of integrating different modalities of source code used in our work. RQ3 aims to investigate how our graph optimizing mechanism improves the training results. RQ4 is to test and verify the robustness of our proposed model when varying the comment length, code length, VFG nodes number and the longest path length of VFG. RQ5 aims to evaluate the performance of DEGRAPHCS compared with the state-of-the-art models in manual evaluation.

4.1 Experimental Setup

Here, we first describe our experimental dataset and present three widely used evaluation metrics. Then we describe the details of the implementation and introduce the baseline models for comparison.

4.1.1 Data Collection

As described in Section 3, our DEGRAPHCS model needs a large-scale training corpus which contains sufficient code fragments and the corresponding comments descriptions. Unfortunately, we cannot get access to the datasets collected by the existing works like MMAN and UNIF. In fact, we have also considered the dataset released by DeepCS, while this dataset only contains the cleaned Java code. It is hard for us to generate the LLVM IR without raw data. Therefore, we re-construct a corpus of C code snippet, which are crawled from GitHub, to verify the performance of our proposed model. We choose C code snippets because C language is popular and LLVM IR has been widely used on C language.

To build the required dataset, we collect high-star C projects from GitHub (a popular open source projects hosting platform). Then we collect our dataset by selecting the C methods which contain the corresponding comment descriptions and can be compiled into LLVM IR from projects. For each C method c , we treat the first sentence appeared in the comment as the corresponding natural language query q since it typically describes the functionality implemented by the method [23, 24]. To reduce bad comments as much as possible, we use regular expression to delete comments which are not related to method function. Furthermore, we filter the (q, c) pairs by the following rules:

- (q, c) will be filtered out if the code snippet c is a constructor or a test method.
- (q, c) will be filtered out if the length of the query q is less than 3 words or longer than 30 words.
- (q, c) will be filtered out if the length of the code snippet c is less than 5 lines or longer than 30 lines.
- If a (q, c) pair appears multiple times in the dataset, we will remove the duplication.

After collecting the corpus of commented code snippets, we then extract LLVM IR for our proposed model and other features of the source code for the baseline models, *i.e.*, method name, tokens, AST and S-CFG. Finally, we obtain 41,152 samples which are more than the 28527 samples in MMAN[5]. And since the 41,152 samples are distributed in 1554 open source projects, it is general to train our model. Following [5], we shuffle our dataset and split it into training set with 39,152 pairs and test set with 2,000 pairs respectively. In order to avoid the bias resulting from evaluating on isolated dataset, we resort to the automatic evaluation metrics on corpus with the ground truth.

Furthermore, for performing manual evaluation, we first randomly select 100 descriptions from the test set, then we carefully choose 50 descriptions which are easiest to understand. And we rewrite the 50 descriptions (*e.g.*, add some conjunction) as test queries to ensure that they are similar enough to the real-world user queries. We construct a search codebase containing 30799 C code snippets without the training samples to guarantee the fairness. 5 experienced participants are required to select the relevant results returned by each model and record the score.

4.1.2 Evaluation Metrics

We choose two common metrics to measure the performance of code search: SuccessRate@ k , and Mean Reciprocal Rank (MRR).

For the automatic and manual evaluation, we adopt both SuccessRate@ k and MRR to assess the performance of a code search model with respect to a set of queries. SuccessRate@ k represents the percentage of queries for which more than one correct snippet succeed to exist in the top k ranked snippets returned by a search model, which is calculated as: SuccessRate @ k = $\left(\frac{1}{|Q|} \sum_{q=1}^{|Q|} \delta(\text{Rank}_q \leq k) \right)$, where Q denotes the set of queries in our automatic evaluation; Rank_q denotes the highest rank of the hit snippets in the returned snippet list for the query; $\delta()$ denotes a indicator function that returns 1 if the Rank of the q_{th} query (Rank_q) is smaller than k otherwise returns 0. SuccessRate@ k is important because a better code search engine should allow developers to find the desired snippet by inspecting fewer results. Following MMAN, we evaluate SuccessRate@1, SuccessRate@5, SuccessRate@10 respectively and a higher SuccessRate@ k value implies a better performance of the code search model.

We also use MRR to measure the ranking of the search results of each model. MRR is the average of the reciprocal ranks of all queries Q . The reciprocal ranks is the inverse of the highest rank of hit code, *i.e.*, Rank . The computation of MRR is: $\text{MRR} = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \frac{1}{\text{Rank}_q}$, where Q denotes the set of queries in the automatic evaluation; Rank_q denotes the rank of the ground-truth code corresponded to the the q_{th} query. The higher the MRR value, the better the code search performance is.

4.1.3 Baseline Models

We compare the effectiveness of DEGRAPHCS with 3 state-of-the-art deep learning-based code search methods:

- **DeepCS**: DeepCS [2] is a deep code search engine using deep neural networks. Instead of matching text similarities like traditional works, DeepCS learns a unified vector representation of both source code and the corresponding natural language query. We use the official implementation provided by the authors¹.
- **UNIF**: UNIF [3] is a supervised extension of the base NCS technique [25]. UNIF maintains significantly lower complexity than previous sequence-of-words-based networks by using a bag-of-words-based network.
- **MMAN**: MMAN is a novel multi-modal neural network for code search [5]. MMAN proposes an attention mechanism to incorporate multiple features including code tokens, AST and S-CFG to learn a more comprehensive representation for code understanding.

Moreover, to answer RQ2 (how effective is our graph-based integration compared with multi-modal attention) and RQ3 (how graph optimization in DEGRAPHCS affects its effectiveness), we compare DEGRAPHCS with some of its variants as follows:

- **MMAN (Token+V-DFG+V-CFG)**: In this variant, we exploit the features of DEGRAPHCS, *i.e.*, tokens of code, variable-based data and control flow graph and fuses them in a multi-modal neural network with attention mechanism. In other words, the features used in MMAN are replaced with the features used in DEGRAPHCS. This variant is used to validate that our graph building and optimizing mechanism is more effective than the previous multi-modal incorporation mechanism.
- **DEGRAPHCS-noGO**: In this variant, we remove the graph optimization mechanism from DEGRAPHCS. This variant is used to validate that it is necessary to remove the redundant information in the initial constructed graph and verify how this mechanism can improve the training results.

4.1.4 Implementation Details

To train our proposed model, we first shuffle the training data and set the mini-batch size to 16. We build two separate vocabularies for comments and LLVM IR tokens and limit their size of vocabulary to 10,000 and 15,000 respectively to store the most frequently appeared tokens in the training dataset. For each batch, comments are padded with a special token “PAD” to the maximum length which is set to 30 in our experiments. All tokens in our dataset are converted to lower case and parsed into a sequence of tokens according to camel case and “_” if exists. We set the word embedding size to 300. For LSTM and GGNN unit, we set the hidden size to 512. Besides, we set 5 rounds of iteration for GGNN. The margin is set to 0.6. We update the parameters via AdamW optimizer [26] with the learning rate 0.0003. All the models in this paper are trained for 200 epochs. All the experiments are implemented using Pytorch 1.6 framework with Python 3.6, and the experiments are conducted on a server with one Nvidia Tesla V100 GPU, running on Ubuntu 18.04.

4.2 Experimental Results

4.2.1 RQ1: Comparison with Baselines

RQ1 aims to investigate whether DEGRAPHCS outperforms the state-of-the-art deep code search models. We evaluate DEGRAPHCS on the test set, which consists of 2,000 pairs of code snippets and the corresponding descriptions. In this automatic evaluation, we treat each description as an input query, and its corresponding code snippet as the ground truth. We report our evaluation results in Table I. The columns R@1, R@5, R@10 show the values of SuccessRate@k over all queries when k is set to 1, 5 and 10, separately. The column MRR presents the MRR value of each model. From Table I, we can clearly find that DEGRAPHCS beats existing code search methods to a large extent on all the metrics. In concrete, DEGRAPHCS obtains an MRR of 51.73%, which is much better than that of DeepCS (32.68%), UNIF (41.93%), MMAN (37.97%). As for SuccessRate@k, DEGRAPHCS improves the state-of-the-art R@1 score from 34.05% (obtained by MMAN) to 43.05%. For 43.05%/61.75%/68.10% of the test queries, the relevant code snippets can be found within the top 1/5/10 returned results by DEGRAPHCS. The results confirm that the improvement achieved by DEGRAPHCS is statistically and consistently significant, which indicates the effectiveness of DEGRAPHCS.

¹<https://github.com/guxd/deep-code-search>

Table 1: Comparison of the overall performance between our model and baselines on automatic evaluation metrics

Method	R@1	R@5	R@10	MRR
DeepCS	0.2350	0.4185	0.5045	0.3268
UNIF	0.3250	0.5175	0.5980	0.4193
MMAN	0.3405	0.5325	0.6130	0.4342
deGraphCS	0.4305	0.6175	0.6810	0.5173

Table 2: Effect of Graph Integration

Method	R@1	R@5	R@10	MRR
MMAN(Token+V-DFG+V-CFG)	0.3380	0.5250	0.5990	0.4277
deGraphCS	0.4305	0.6175	0.6810	0.5173

4.2.2 RQ2: Effects of Integration

To answer this question, we perform experiments over different combination mechanism of the features (*i.e.*, code tokens, variable-based data and control graph) by two models: DEGRAPHCS and MMAN (Token+V-DFG+V-CFG). Table II shows the overall performance of the two approaches.

From Table II, we can observe that the MRR of MMAN (Token+V-DFG+V-CFG) decreases by 8.96% compared with DEGRAPHCS. In terms of SuccessRate@k, MMAN (Token+V-DFG+V-CFG) achieves a SuccessRate@1/5/10 of 33.80%, 52.50% and 59.90% respectively, much lower than those of DEGRAPHCS. It means that more relevant code snippets can be returned by DEGRAPHCS. Therefore, integrating the three features into one graph preforms better than roughly fusing them by a single attention layer.

4.2.3 RQ3: Effect of Graph Optimization Component

To demonstrate the effectiveness of the flow graph optimizing mechanism constructed from LLVM IR, we perfrom experiments by comparing DEGRAPHCS with the version removing the graph optimization mechanism DEGRAPHCS-noGO. We present the overall comparison results in Table III.

From Table III, we can see that DEGRAPHCS-noGO achieves average success rate of 36.20%, 55.85%, 62.10% when the top 1, 5, 10 results are inspected, respectively, while DEGRAPHCS achieves 18.92%, 10.56% and 9.66% improvement in terms of SuccessRate@1, SuccessRate@5 and SuccessRate@10, respectively. In terms of MRR, DEGRAPHCS achieves a 13.77% improvement compared with DEGRAPHCS-noGO.

The results demonstrate that removing redundant information and optimizing node contents of the initial constructed flow graph will make our proposed model focus more on the useful and fine-grained correlations between source code and the comment descriptions. Besides, the performance of DEGRAPHCS-noGO still outperforms state-of-the-art models, which further verifies the robustness of our proposed code representation method. Besides, we have also made some statistics and found that the total number of nodes in graph has been reduced by 51.88%, which results in nearly half decreasing of training time.

4.2.4 RQ4: Model Robustness

To analyze the sensitivity of DEGRAPHCS, we explore four parameters (*i.e.*, comment length, code length, number of VFG nodes, the length of the longest path in VFG) which may have an impact on the code and comment representation. Fig. 5 illustrate the performance of DEGRAPHCS on the basis of different evaluation metrics with varying parameters. From Fig. 5(a)(b)(c) we can find that in most cases, DEGRAPHCS maintains a stable performance even though the coment length, code length or node number increases dramatically, which can be attributed to the superiority of our variable-based flow graph. The length of longest path between any two nodes is used to measure the complexity of our flow graph. From Fig. 5(d), we can see that with the length of the longest path increases (*i.e.*, the difficulty of the graph embedding increases), the performance of DEGRAPHCS decreases slightly but overall maintains a stable level. Overall, the results in this subsection further verify the robustness of our proposed model.

Table 3: Effect of Graph Optimization

Method	R@1	R@5	R@10	MRR
deGraphCS-noGO	0.3620	0.5585	0.6210	0.4547
deGraphCS	0.4305	0.6175	0.6810	0.5173

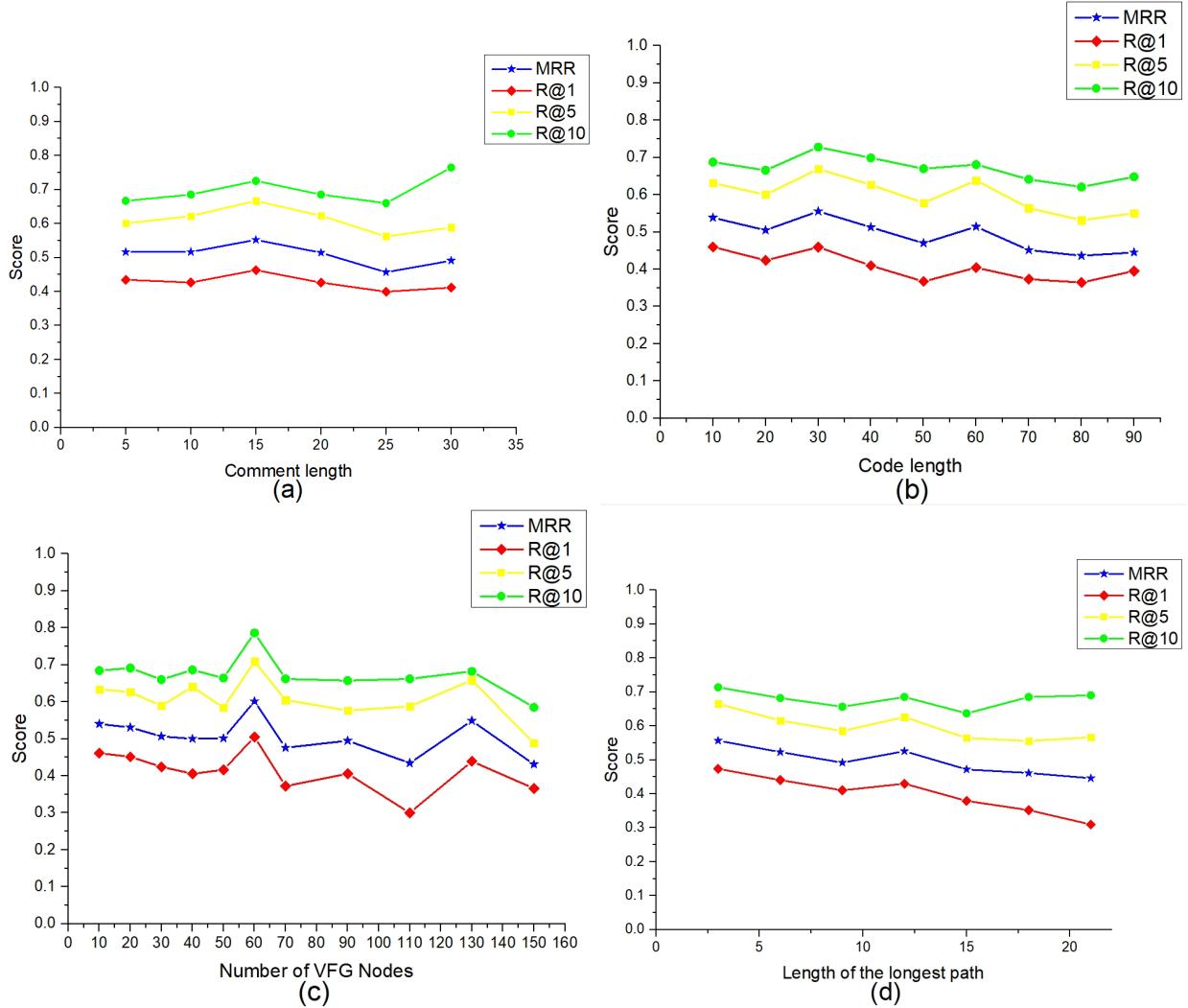


Figure 5: Experimental results of deGraphCS on different metrics w.r.t. varying number of nodes and code lengths.

4.2.5 RQ5: Human Evaluation

DEGRAPHCS has shown great utility in the aforementioned automatic evaluation experiments. However, in reality, the questions of how useful a returned code snippet is are likely best answered by human programmers. Thus, we conduct a user study where 5 experienced programmers are asked to grade the utility of code fragments returned by the baseline methods, *i.e.*, DeepCS, UNIF, MMAN and our proposed model DEGRAPHCS.

In concrete, we first build a search platform which includes a search codebase of 30,799 C language functions and 50 queries select from the test set as benchmark. By using different models, the platform will recommend top 10 code snippets to users according to the query. Then, we recruit 5 graduate students from our school, who have rich experience on C projects and are competent enough to perform the user study, to evaluate the effectiveness of these models.

Table 4: Comparison of the overall performance between our model and baselines on manual evaluation (SuccessRate@10|MRR)

Method	P1	P2	P3	P4	P5	Avg.
DeepCS	0.40 0.21	0.34 0.18	0.34 0.17	0.52 0.36	0.48 0.35	0.42 0.26
UNIF	0.52 0.36	0.52 0.27	0.48 0.29	0.60 0.39	0.58 0.39	0.54 0.34
MMAN	0.58 0.41	0.52 0.33	0.48 0.37	0.64 0.43	0.64 0.44	0.57 0.40
deGraphCS	0.66 0.47	0.62 0.46	0.56 0.36	0.70 0.51	0.70 0.61	0.65 0.48

To simulate real scenarios, we let 5 participants (denoted as P1 - P5) use the four models in turn, to search related code of 50 queries. All participants have over 2-years/5-projects C programming experience. During the evaluation, we make sure that participants did not know which model the results searched by. For each query, they need to inspect the top 10 results returned by each model and label those results they believe are relevant to the query. Table IV shows the overall performance achieved by each model in terms of SuccessRate@10 and MRR.

From Table IV, we can draw the following conclusions: (a) under the experimental setting, DEGRAPHCS answer more user queries with an average SuccessRate@10 of 0.65 and the improvements compared with MMAN, UNIF and DeepCS is 14%, 20% and 55%, respectively; (b) DEGRAPHCS achieves a better code search performance with an average MRR of 0.48. In conclusion, our proposed model maintains higher practice value in simulated code search scenario.

We further choose several examples to illustrate the superiority of DEGRAPHCS on the searching results. Fig. 6 shows the searched results of DEGRAPHCS, DeepCS, UNIF and MMAN on the query “allocate memory for the file descriptors”. We can see that the result returned by DEGRAPHCS is exactly what the users need. However, the results returned by DeepCS, UNIF and MMAN do not realize the queried function and only focus on the shallow information (*i.e.*, keyword “fd”, which is related to file descriptor). In concrete, the core functionality related to the keywords in the query is squared in red. The baseline methods can only retrieve the functions that realize the file creation, open or memory allocation for other data structures instead of file descriptor. Fig. 7 shows the rank 1 and rank 2 searched results of DEGRAPHCS on the query “calculate checksum of checkpoint”. We can see that both the two retrieved code snippets realize the function of buffer checking and computing. However, compared with the rank 1 result, whose tokens are well-named (*i.e.*, checksum, buffer) and obviously matched with the keywords in the query, the variables in rank 2 code snippet are much more obscure (*i.e.*, cksum, b). Since the tokens in rank 2 code are not named following natural languages rules, it is hard for the models to utilize the token information for matching. The fact that DEGRAPHCS can retrieve the obscure code snippet related to the query on the semantics demonstrates that the data and control flow features are essential in code search and DEGRAPHCS can fully exploit the useful information in our variable-based flow graph.

4.3 Threats to Validity

DEGRAPHCS may suffer from three threats to validity. The first one lies in the scalability of our proposed approach. The LLVM IR can only be extracted from a whole program with complete dependencies. Therefore, it is difficult to extend our precise code representation model to some places where LLVM IR cannot be successfully extracted, such as many single code snippets in Stack Overflow. In fact, the difficulty can be solved by automatically adding the interface for those missing class objects and methods. And we plan to overcome these compilation problem to generate more dataset in the near future. The second threat is that DEGRAPHCS is currently trained and tested only on C programs. However, our work is based on LLVM IR of the code, which is independent of the source programming language. Thus, our code representation method can be easily transferred to other languages like Python and Java *etc.* For example, to transfer our work to Java language, we can use Soot[27] (a Java optimization framework) to generate jimple which is intermediate representation of Java. And to transfer to python language, we can analysis and utilize the bytecode of python. We can build variable-based flow graph based on the above intermediate representation. The third threat lies in the model generalization ability. We construct a test dataset only consisting of 2,000 C code snippets, which may not be enough to represent most programming tasks. We plan to extend the dataset in the near future.

```

alloc_fds(int **fds, int n_fds){
    if (n_fds) {
        *fds = (int *)malloc(sizeof(int) * n_fds);
        if (*fds == NULL) return OP_ERROR;
    }
    return OP_SUCCESS;
}

```

(a) deGraphCS

```

void cmd_mkexec(frommark, tomark, cmd, bang, extra){
    int fd;
    .....
    fd = creat(extra, FILEPERMS);
    if (fd < 0){
        msg("Couldn't create a new \"%s\" file", extra); return;
    savekeys(fd); saveopts(fd);
    .....
    close(fd); msg("Created a new \"%s\" file", extra);
}

```

(b) DeepCS

```

struct bmp_decdta *bmp_alloc(void){
    struct bmp_decdta *bmp = malloc_tmphigh(sizeof(*bmp)); return bmp;
}

```

(c) UNIF

```

int ffwopen(const char *fn, struct buffer *bp){
    int fd; mode_t fmode = DEFFILEMODE;
    .....
    fd = open(fn, O_RDWR | O_CREAT | O_TRUNC, fmode);
    if (fd == -1) {
        ffp = NULL; ewprintf("Cannot open file for writing : %s", strerror(errno));
        return (FIOERR);
    }
    .....
    return (FIOSUC);
}

```

(d) MMAN

Figure 6: An illustrative example shows the comparison between the searched results of (a) deGraphCS, (b) DeepCS, (c) UNIF and (d) MMAN on the query “allocate memory for the file descriptors”.

```

checksum(unsigned short *buffer, int size)
{
    unsigned Long cksum=0;
    while (size > 1) {
        cksum += *buffer++;
        size -= sizeof(unsigned short);
    }
    if (size) {
        cksum += *(unsigned char *)buffer;
    }
    cksum = (cksum >> 16) + (cksum & 0xffff);
    cksum += (cksum >>16);
    return (unsigned short)(~cksum);
}

```

(a) deGraphCS Rank 1

```

uint16_t cksum( char *b, int len ){
    uint16_t sum = 0;
    uint16_t t;
    char *e = b + len;
    b[len] = 0;
    while(b < e){
        t = (b[0] << 8) + b[1];
        sum += t;
        if(sum < t)
            sum++;
        b += 2;
    }
    return ~sum;
}

```

(b) deGraphCS Rank 2

Figure 7: The rank 1 and rank 2 searched results of deGraphCS on the query “calculate checksum of checkpoint”.

5 Related Work

5.1 Code Search

As vast repositories of open source code are available, a lot of works are proposed to search code to help developers. There are many works that try to search code according to user query. The traditional code search methods are mainly based on information retrieval and natural language processing technologies [8, 9, 10, 11, 12, 13, 14, 28]. The methods mainly take text and structure characteristics of source code into consideration. [28] proposes a code search engine Sourcerer which extracts fine-grained structural information from source code. [12] proposes CodeHow, a code search technique which reveals APIs related to user queries according to text similarity, and then applies an extended boolean model to utilize API information in code searching. NCS[25] utilizes natural language processing technologies to embed the code and query into vector, and then search the code snippet by comparing similarity of the vector representations. The above information retrieval-based methods treat both source code and query as natural language. It is hard to deeply understand the semantic information of source code.

Since traditional code search based on syntactic are easy to return vague or irrelevant result, many works[29, 30, 31, 32, 33] are proposed for semantic code search which is depends on specifications. For example, to search semantic related code, [29] proposed a architecture for semantics-based code search. The architecture utilize many different specifications which includes keywords, method signatures and test cases *etc*. The semantic search based on specifications performs well for finding relevant code but requires developers to write complex specifications. To reduce the requirement of specifications, [30] proposes a new approach in which programmers only are required to specify lightweight, incomplete specifications which are in the form of input/output pairs and/or partial program fragment. Then the approach uses an SMT solver to automatically identify programs. While it also requires extra specifications to understand code semantics. Compared with the prior semantic code search, our method uses a static code-level analysis on source code, without the need to input extra specifications or run the code snippets. And our work focus on the scenario that users can simply use natural language to describe their intents. Last, we obtain code semantic by using a deep learning model to learn a representation of code instead of depending on the input and analysis of specifications.

Plenty of works are proposed to enrich information of the queries by refining including query expansion and reformulation[34, 35, 36, 37, 31, 38, 39]. In concrete, [34] trains a recommender (Refoqus) based on machine learning technologies, Refoqus can recommend a reformulation strategy according to the properties of queries. [37] utilizes synonyms generated from WordNet to extend the queries. [31] utilizes the feedback of users to reformulate queries.

Recently, to understand the deep semantics of the code and query, deep learning technologies have been applied to code search [1, 40, 2, 41, 42, 43, 5, 4, 25, 3]. [1] proposes BVAE including two Variational AutoEncoders (VAEs) to model source code and natural language respectively, and jointly train two model to capture the closeness between the latent variables of the code and description. Many works measure semantic similarity of source code and query through joint embedding and deep learning technologies. CODEnn [2] extracts tokens, filenames and API sequences of code as the features, and embed these information and queries into a shared space so that code snippet can be retrieved by the vector of query. UNIF [3] extends NCS and further fine-tuned embedding of code and query by jointly deep learning. To utilize more information of source code, [5] proposed MMAN which use a multi-modal (tokens,AST and CFG) to represent source code. Similarly, we embed the source code and query into a common space to mine the semantic relationship. However, these works can not precisely represent the semantics of source code.

5.2 Code Representation

Deep learning techniques on program analysis have attracted increasing attention. Many works focus on the representation of source code to perform a further software engineering research [44, 45, 46, 47, 48, 49, 50, 51, 52, 53]. [52] adopts a RNN and n-gram model for code completion. To capture the structure information of code, [48] proposes a novel novel tree-based convolutional neural network (TBCNN) to represent the ASTs of source code. [51] proposes a framework (CDLH) which incorporates an AST-based LSTM to exploit the lexical and syntactical information. More recently, to extract more information, [5] proposes MMAN to combine multiple semantic information of code, which includes tokens, AST, and CFG of the source code. These researches focus on different representations of source code to capture the structural, syntactical and semantic information. However, these works can not represent the code precisely on semantics. Therefore, we pay more attention to the semantics of code, and propose our variable-based flow graph method.

Several recent work[54, 55, 56, 57] have tried to utilize intermediate representation to represent code. [54] construct dependency graph to represent code based on instructions of Java bytecode. They use graph to record data and control dependencies between instructions, so that they can detect the similarity of code by use a subgraph isomorphism

algorithm to analyze the similarity of dependency graph. [55] also focuses on building a graph based on LLVM IR, to represent code. And they obtain well performance when they use skip-gram model to learn the graph representation and apply to code task. Different with these works, we aim to make a more precise code search, thus we construct a different graph based on variable instead of instructions, which ensures that the granularity of both comments and code representations in code search is consistent. And compared with the prior used method, we use a deep learning model (GGNN) to learn the semantic representation of code. Furthermore, to obtain a more precise representation of code, we make an optimization on the graph to decrease the noises and improve the training efficiency.

6 Conclusion and Future Work

In this paper, we propose a deep graph neural network named DEGRAPHCS for code search. Instead of considering structural features of source code such as AST, DEGRAPHCS incorporates multiple semantic features, *i.e.*, tokens, variable-based data and control flow extracted from LLVM IR of code into flow graph. Furthermore, we put forward an optimization to remove the redundant information of the graph, followed by a gated graph neural network with attention mechanism to capture the critical information of code. In addition, we use a unified framework to learn the representation of natural language query and corresponding code snippet. We conduct several experiments on trained models, and the results of automatic evaluation and manual evaluation both demonstrate that our proposed approach is effective and outperforms the state-of-the-art approaches.

In future, we plan to investigate the performance of the DEGRAPHCS on other dataset of different programming languages, *e.g.*, Python and Java. We also plan to extend the variable-based flow graph we designed to solve other software engineering problems, *e.g.*, code translation[55], API recommendation [43, 58], and code clone detection [46].

References

- [1] Qingying Chen and Minghui Zhou. A neural framework for retrieval and summarization of source code. In *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 826–831. IEEE, 2018.
- [2] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. Deep code search. In *2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE)*, pages 933–944. IEEE, 2018.
- [3] Jose Cambronero, Hongyu Li, Seohyun Kim, Koushik Sen, and Satish Chandra. When deep learning met code search. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 964–974, 2019.
- [4] Jianhang Shuai, Ling Xu, Chao Liu, Meng Yan, Xin Xia, and Yan Lei. Improving code search with co-attentive representation learning. In *28th International Conference on Program Comprehension (ICPC)*, 2020.
- [5] Yao Wan, Jingdong Shu, Yulei Sui, Guandong Xu, Zhou Zhao, Jian Wu, and Philip Yu. Multi-modal attention network learning for semantic source code retrieval. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 13–25. IEEE, 2019.
- [6] Joel Brandt, Mira Dontcheva, Marcos Weskamp, and Scott R Klemmer. Example-centric programming: integrating web search into the development environment. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 513–522, 2010.
- [7] Brock Angus Campbell and Christoph Treude. Nlp2code: Code snippet content assist via natural language tasks. In *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 628–632. IEEE, 2017.
- [8] Wing-Kwan Chan, Hong Cheng, and David Lo. Searching connected api subgraph via text phrases. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, pages 1–11, 2012.
- [9] Reid Holmes, Rylan Cottrell, Robert J Walker, and Jorg Denzinger. The end-to-end use of source code examples: An exploratory study. In *2009 IEEE International Conference on Software Maintenance*, pages 555–558. IEEE, 2009.
- [10] Iman Keivanloo, Juergen Rilling, and Ying Zou. Spotting working code examples. In *Proceedings of the 36th International Conference on Software Engineering*, pages 664–675, 2014.
- [11] Xuan Li, Zerui Wang, Qianxiang Wang, Shoumeng Yan, Tao Xie, and Hong Mei. Relationship-aware code search for javascript frameworks. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 690–701, 2016.

- [12] Fei Lv, Hongyu Zhang, Jian-guang Lou, Shaowei Wang, Dongmei Zhang, and Jianjun Zhao. Codehow: Effective code search based on api understanding and extended boolean model (e). In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 260–270. IEEE, 2015.
- [13] Collin McMillan, Mark Grechanik, Denys Poshyvanyk, Chen Fu, and Qing Xie. Exemplar: A source code search engine for finding highly relevant applications. *IEEE Transactions on Software Engineering*, 38(5):1069–1087, 2011.
- [14] Collin McMillan, Mark Grechanik, Denys Poshyvanyk, Qing Xie, and Chen Fu. Portfolio: finding relevant functions and their usage. In *Proceedings of the 33rd International Conference on Software Engineering*, pages 111–120, 2011.
- [15] Chris Lattner and Vikram Adve. LLVM: A compilation framework for lifelong program analysis & transformation. In *International Symposium on Code Generation and Optimization, 2004. CGO 2004.*, pages 75–86. IEEE, 2004.
- [16] Tal Ben-Nun, Alice Shoshana Jakobovits, and Torsten Hoefer. Neural code comprehension: A learnable representation of code semantics. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 3585–3597. Curran Associates, Inc., 2018.
- [17] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. *CoRR*, abs/1503.00075, 2015.
- [18] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*, 2015.
- [19] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Trans. Program. Lang. Syst.*, 13(4):451–490, October 1991.
- [20] Junyoung Chung, Çağlar Gülcabay, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
- [21] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [22] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [23] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. Codesearchnet challenge: Evaluating the state of semantic code search. *CoRR*, abs/1909.09436, 2019.
- [24] Xing Hu, Ge Li, Xin Xia, David Lo, and Zhi Jin. Deep code comment generation. In *2018 IEEE/ACM 26th International Conference on Program Comprehension (ICPC)*, pages 200–20010. IEEE, 2018.
- [25] Saksham Sachdev, Hongyu Li, Sifei Luan, Seohyun Kim, Koushik Sen, and Satish Chandra. Retrieval on source code: a neural code search. In *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 31–41, 2018.
- [26] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017.
- [27] Soot. <https://github.com/soot-oss/soot>.
- [28] Sushil Bajracharya, Trung Ngo, Erik Linstead, Yimeng Dou, Paul Rigor, Pierre Baldi, and Cristina Lopes. Sourcerer: a search engine for open source code supporting structure-based search. In *Companion to the 21st ACM SIGPLAN symposium on Object-oriented programming systems, languages, and applications*, pages 681–682, 2006.
- [29] S. P. Reiss. Semantics-based code search. In *2009 IEEE 31st International Conference on Software Engineering*, pages 243–253, 2009.
- [30] K. T. Stolee. Finding suitable programs: Semantic search with incomplete and lightweight specifications. In *2012 34th International Conference on Software Engineering (ICSE)*, pages 1571–1574, 2012.
- [31] Shaowei Wang, David Lo, and Lingxiao Jiang. Active code search: incorporating user feedback to improve code search relevance. In *Proceedings of the 29th ACM/IEEE international conference on Automated software engineering*, pages 677–682, 2014.
- [32] Otávio Augusto Lazzarini Lemos, Sushil Bajracharya, Joel Ossher, Paulo Cesar Masiero, and Cristina Lopes. A test-driven approach to code search and its application to the reuse of auxiliary functionality. *Information and Software Technology*, 53(4):294 – 306, 2011.

- [33] Otávio Augusto Lazzarini Lemos, Sushil Krishna Bajracharya, and Joel Ossher. Codegenie: : a tool for test-driven source code search. In Richard P. Gabriel, David F. Bacon, Cristina Videira Lopes, and Guy L. Steele Jr., editors, *Companion to the 22nd Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2007, October 21-25, 2007, Montreal, Quebec, Canada*, pages 917–918. ACM, 2007.
- [34] Sonia Haiduc, Gabriele Bavota, Andrian Marcus, Rocco Oliveto, Andrea De Lucia, and Tim Menzies. Automatic query reformulations for text retrieval in software engineering. In *2013 35th International Conference on Software Engineering (ICSE)*, pages 842–851. IEEE, 2013.
- [35] Emily Hill, Lori Pollock, and K Vijay-Shanker. Improving source code search with natural language phrasal representations of method signatures. In *2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011)*, pages 524–527. IEEE, 2011.
- [36] Emily Hill, Manuel Roldan-Vega, Jerry Alan Fails, and Greg Mallet. NL-based query refinement and contextualized code search results: A user study. In *2014 Software Evolution Week-IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering (CSMR-WCRE)*, pages 34–43. IEEE, 2014.
- [37] Meili Lu, Xiaobing Sun, Shaowei Wang, David Lo, and Yucong Duan. Query expansion via wordnet for effective code search. In *2015 IEEE 22nd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*, pages 545–549. IEEE, 2015.
- [38] Timothy Dietrich, Jane Cleland-Huang, and Yonghee Shin. Learning effective query transformations for enhanced requirements trace retrieval. In *2013 28th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 586–591. IEEE, 2013.
- [39] Otávio AL Lemos, Adriano C de Paula, Felipe C Zanichelli, and Cristina V Lopes. Thesaurus-based automatic query expansion for interface-driven code search. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, pages 212–221, 2014.
- [40] Daniel DeFreez, Aditya V Thakur, and Cindy Rubio-González. Path-based function embedding and its application to specification mining. *arXiv preprint arXiv:1802.07779*, 2018.
- [41] Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. Deep api learning. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 631–642, 2016.
- [42] Jordan Henkel, Shuvendu K Lahiri, Ben Liblit, and Thomas Reps. Code vectors: understanding programs through embedded abstracted symbolic traces. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 163–174, 2018.
- [43] Qiao Huang, Xin Xia, Zhenchang Xing, David Lo, and Xinyu Wang. Api method recommendation without worrying about the task-api knowledge gap. In *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 293–304. IEEE, 2018.
- [44] Hamel Husain and Ho-Hsiang Wu. How to create natural language semantic search for arbitrary objects with deep learning. *Retrieved November, 5:2019*, 2018.
- [45] Hamel Husain and Ho-Hsiang Wu. Towards natural language semantic code search. *Retrieved November, 5:2019*, 2018.
- [46] Martin White, Michele Tufano, Christopher Vendome, and Denys Poshyvanyk. Deep learning code fragments for code clone detection. In *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 87–98. IEEE, 2016.
- [47] Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S Yu. Improving automatic source code summarization via deep reinforcement learning. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, pages 397–407, 2018.
- [48] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. Convolutional neural networks over tree structures for programming language processing. *arXiv preprint arXiv:1409.5718*, 2014.
- [49] Chris Piech, Jonathan Huang, Andy Nguyen, Mike Phulsuksombati, Mehran Sahami, and Leonidas Guibas. Learning program embeddings to propagate feedback on student code. *arXiv preprint arXiv:1505.05969*, 2015.
- [50] Hoa Khanh Dam, Truyen Tran, and Trang Pham. A deep language model for software code. *arXiv preprint arXiv:1608.02715*, 2016.
- [51] Huihui Wei and Ming Li. Supervised deep features for software functional clone detection by exploiting lexical and syntactical information in source code. In *IJCAI*, pages 3034–3040, 2017.

- [52] Veselin Raychev, Martin Vechev, and Eran Yahav. Code completion with statistical language models. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 419–428, 2014.
- [53] Michele Tufano, Cody Watson, Gabriele Bavota, Massimiliano Di Penta, Martin White, and Denys Poshyvanyk. Deep learning similarities from different representations of source code. In *2018 IEEE/ACM 15th International Conference on Mining Software Repositories (MSR)*, pages 542–553. IEEE, 2018.
- [54] Fang-Hsiang Su, Jonathan Bell, Kenneth Harvey, Simha Sethumadhavan, Gail Kaiser, and Tony Jebara. Code relatives: Detecting similarly behaving software. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, FSE 2016, page 702–714, New York, NY, USA, 2016. Association for Computing Machinery.
- [55] Xinyun Chen, Chang Liu, and Dawn Song. Tree-to-tree neural networks for program translation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 2547–2557. Curran Associates, Inc., 2018.
- [56] A. Schäfer, W. Amme, and T. S. Heinze. Detection of similar functions through the use of dominator information. In *2020 IEEE International Conference on Autonomic Computing and Self-Organizing Systems Companion (ACSOS-C)*, pages 206–211, 2020.
- [57] P. M. Caldeira, K. Sakamoto, H. Washizaki, Y. Fukazawa, and T. Shimada. Improving syntactical clone detection methods through the use of an intermediate representation. In *2020 IEEE 14th International Workshop on Software Clones (IWSC)*, pages 8–14, 2020.
- [58] Liang Cai, Haoye Wang, Qiao Huang, Xin Xia, Zhenchang Xing, and David Lo. Biker: A tool for bi-information source based api method recommendation. ESEC/FSE 2019, New York, NY, USA, 2019. Association for Computing Machinery.