



# Fast code recommendation via approximate sub-tree matching\*

Yichao SHAO<sup>†1,2,3</sup>, Zhiqiu HUANG<sup>†‡1,2,3</sup>, Weiwei LI<sup>1,2,3</sup>, Yaoshen YU<sup>1,2,3</sup>

<sup>1</sup>School of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing 211100, China

<sup>2</sup>Key Laboratory of Safety-Critical Software, Ministry of Industry and Information Technology, Nanjing 211100, China

<sup>3</sup>Collaborative Innovation Center of Novel Software Technology and Industrialization, Nanjing 210016, China

<sup>†</sup>E-mail: shaoyichao@nuaa.edu.cn; zqhuang@nuaa.edu.cn

Received Aug. 7, 2021; Revision accepted Mar. 24, 2022; Crosschecked July 15, 2022

**Abstract:** Software developers often write code that has similar functionality to existing code segments. A code recommendation tool that helps developers reuse these code fragments can significantly improve their efficiency. Several methods have been proposed in recent years. Some use sequence matching algorithms to find the related recommendations. Most of these methods are time-consuming and can leverage only low-level textual information from code. Others extract features from code and obtain similarity using numerical feature vectors. However, the similarity of feature vectors is often not equivalent to the original code's similarity. Structural information is lost during the process of transforming abstract syntax trees into vectors. We propose an approximate sub-tree matching based method to solve this problem. Unlike existing tree-based approaches that match feature vectors, it retains the tree structure of the query code in the matching process to find code fragments that best match the current query. It uses a fast approximation sub-tree matching algorithm by transforming the sub-tree matching problem into the match between the tree and the list. In this way, the structural information can be used for code recommendation tasks that have high time requirements. We have constructed several real-world code databases covering different languages and granularities to evaluate the effectiveness of our method. The results show that our method outperforms two compared methods, SENSORY and Aroma, in terms of the recall value on all the datasets, and can be applied to large datasets.

**Key words:** Code reuse; Code recommendation; Tree similarity; Structure information

<https://doi.org/10.1631/FITEE.2100379>

**CLC number:** TP311

## 1 Introduction

In software development, developers tend to reuse existing code which can achieve the desired functionality or behavior to assist their development procedure. Research shows that, on average, software developers spend about 19% of their development time on web searches, looking mainly for code examples for their tasks (Rahman et al., 2016). Therefore, an automatic code snippet recommendation tool

could help developers greatly improve development efficiency.

However, searching instructive code based on a user's programming context is not always easy. Most of the current code recommendation tools are based on textual matching. They represent the context code fragment as tokens or lines of code, calculating the code similarity through sequence matching or a bag-of-words model, and return the most relevant code snippets (Ye and Fischer, 2002; Holmes and Murphy, 2005; Antunes et al., 2014; Rahman and Roy, 2014). However, these methods fail to capture high-level structural information, which means that it would be hard to obtain the best result when there is no highly similar match in the code repository.

<sup>‡</sup> Corresponding author

\* Project supported by the National Natural Science Foundation of China (No. 61772270)

ORCID: Yichao SHAO, <https://orcid.org/0000-0002-4553-5602>; Zhiqiu HUANG, <https://orcid.org/0000-0001-6843-1892>

© Zhejiang University Press 2022

Abstract syntax tree (AST) carries the structural information of code, but a traditional tree matching process consumes too much time for code recommendation tasks. To solve this problem, some researchers proposed to extract features from the AST of the code and used these, rather than trees, for matching (Jiang LX et al., 2007; Luan et al., 2018; Ďuračik et al., 2020). However, because the ideal candidate code length should be larger than the query code, these candidate codes often contain more sub-tree structures and feature patterns. That is to say, although two code snippets may obtain a high similarity score, they may not be related code because the distribution of these features in candidate code can be different from that in query code. Therefore, methods using such techniques may bring more false-positive results.

In this study, we propose a novel code recommendation method based on AST sub-tree matching. Compared with methods that transform all ASTs into multiple features, we retain the tree structure of the query code in the matching process. A candidate code with a complete query code structure will gain a higher similarity score than those that share the same feature sets but have different overall structures, and thus will have a higher chance of being recommended. We use hash values to record every sub-tree in the AST and speed up the matching process by comparing the hash value only of sub-trees with the same number of nodes. Our method includes three stages: data processing, coarse-grained search, and fine-grained re-ranking. First, we construct the code database using AST. In this stage, we calculate the hash value and the number of nodes for each sub-tree in AST, and store them using a list structure. This helps reduce the costs of time and space in the matching process. Second, we traverse the AST of a given query code snippet to calculate the similarity between the query code and the code-base's candidate code. In this stage, the similarity measurement of two code snippets is the number of nodes in all the most similar sub-trees of their ASTs. Based on the similarity score, a top- $K$  candidate set will be generated. Third, we obtain the AST preorder traversal sequence for both the query code and each top- $K$  candidate code. We calculate their similarity using the Smith–Waterman (SW) algorithm to

fully mine the sequence and continuity information in the code and re-rank the top- $K$  list. In this way, we can eventually obtain the recommending result list.

We have carried out an extensive empirical evaluation of our method applied to several large datasets, covering different languages and granularities. We compared our method with two strong baseline studies, and the experimental results showed that our method outperforms the two compared methods on all the datasets.

This study has the following main contributions:

1. We introduce a new code recommendation algorithm based on sub-tree hashing and the SW algorithm. It takes programming context as input and recommends relevant code snippets to assist developers in software development.
2. We implement a code recommendation tool for Java and C. Experimental results show that it has good performance in terms of both time consumption and accuracy for different recommending tasks.

## 2 Related works

In this section, we describe some related studies, split into two categories: (1) tree similarity detection; (2) code recommendation.

1. Tree similarity detection. Since tree-like data structure has gained more popularity in recent years, many studies have proposed to focus on similarity detection on trees. However, since most existing methods cannot be extended to large-scale datasets, how to achieve fast tree similarity detection is still an open question. Some researchers use tree edit distances to measure the similarity (Zhang and Shasha, 1989; Shasha et al., 1994; Chen and Zhang, 2014). Such edit-distance-based methods are widely used, but they obtain the similarity through the difference between two trees. This is inappropriate for tasks in which one tree contains another, which happens frequently in code recommendation. Other researchers have proposed to extract features from trees and transform the tree similarity detection problem into feature matching on numerical vectors, which is fast and easy (Jiang LX et al., 2007; Luan et al., 2018; Ďuračik et al., 2020). However, the similarity of feature

sets is not equivalent to the similarity of trees. These methods abandon the original structure of the trees and the continuity information of nodes, and may return false high similarity results.

2. Code recommendation. Code snippet recommendation is a hot issue that attracts many researchers. Most researchers focus on improving the precision and speed of recommending to improve the users' experience. Techniques like information retrieval (Sahavechaphan and Claypool, 2006; Jiang H et al., 2019) and pattern matching (Jiang LX et al., 2007; Āuraćik et al., 2020) are widely used in their work. Most existing methods are based on textual information. They treat code snippets as a set of tokens or code lines (Ye and Fischer, 2002; Holmes and Murphy, 2005; Antunes et al., 2014; Rahman and Roy, 2014; Ai et al., 2019). Some researchers parse the code snippets into AST and use tree-based similarity detection techniques to obtain recommendation results (Luan et al., 2018; Āuraćik et al., 2020). The advantage of such methods is that they exploit the syntax information of the code, but they often consume more time.

### 3 Methods

Fig. 1 illustrates the overall architecture of our framework. To use the codebase efficiently, we first calculate and store the hash value and node number of each sub-tree in their ASTs. Then, we traverse the AST of the given query code to calculate the similarity between the query code and all the candidate code in the codebase. Multiple candidate code

segments with the highest similarity will be selected. Finally, we calculate the AST preorder sequence similarity between the query code and the top- $K$  candidate code using the SW algorithm. We re-rank the candidate code list based on the combination of two kinds of similarity and return the final recommendation result.

Next, we describe the details of each step using the code fragment in Listing 1 as a running example.

**Listing 1** A simple piece of code used as the running example through Section 3

```
for (int i=0; i<10; i++)
    a+=0.5;
end for
```

#### 3.1 Data processing

In this part, we show the procedure for data processing; i.e., how we transform candidate source code files into the form needed for the recommendation.

##### 3.1.1 Tree parsing

Since our method is based on AST, the first thing we need to do is to parse the code. We parse Java code with Javalang (<https://github.com/c2nes/javalang>), and parse C code using Clang (<https://clang.llvm.org/>). Our method can work at either the file level or the method level. For individual Java methods, we add a foo class so that the parser can process them, while for file-level code, we exclude package declarations and import statements, since they are often generated by the integrated development environment (IDE) and not relevant for recommending purposes.

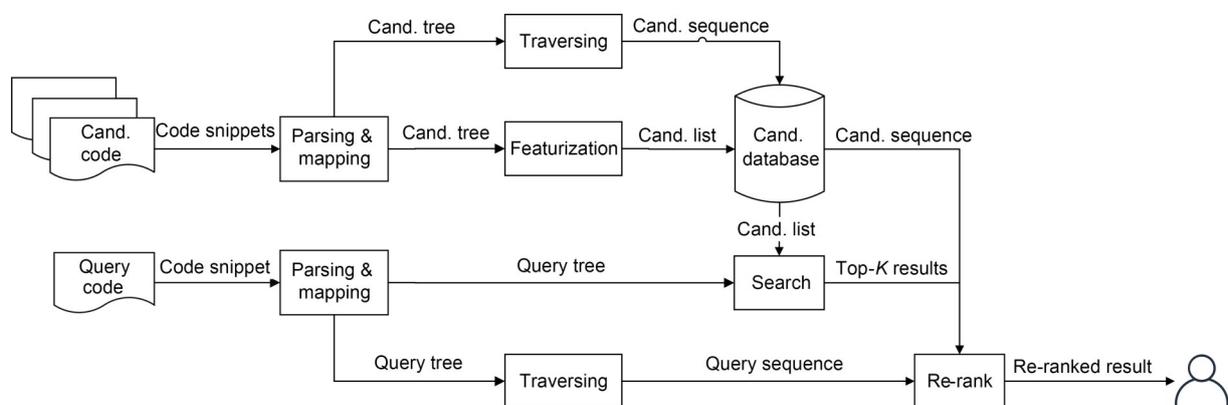


Fig. 1 The framework of code recommendation (Cand.: Candidate)

### 3.1.2 Featurization

We first map each node in the AST to a hash value. Fig. 2 visualizes the simplified parse tree and the corresponding hash value for each node of the code snippet in Listing 1. In principle, different AST nodes should have different values. However, to improve the detection effect, we map similar nodes to the same hash value to avoid the reduction of similarity caused by minor differences among similar nodes (Yang et al., 2018). For example, in Fig. 2, both `FLOATING_LITERAL` and `INTEGER_LITERAL` belong to numerical types, so they are given the same hash value of 13. We classify nodes according to their functions. The specific mapping rules may change according to different languages or parsers. In our experiment, we use Javalang as the parser of Java, and record the classification details in Appendix A.

Nodes under the same classification will be mapped to the same hash value. In a specific implementation, it can be supplemented or adjusted according to the required recommending effect.

Then we calculate the hash value of each sub-tree in the AST. As Fig. 3 shows, we use a `<node number, hash value>` tuple to represent the information of the sub-tree rooted at the current node. We record the node number of each sub-tree and simply define the hash value of a tree as the sum of the hash values of each node in this tree. For the snippet `i<10` in Listing 1, its hash value would be  $\text{Hash}(\text{BINARY\_OPERATOR}) + \text{Hash}(\text{DECL\_REF}) + \text{Hash}(\text{INTEGER\_LITERAL}) = 2 + 17 + 13 = 32$ , and obviously its node number is 3.

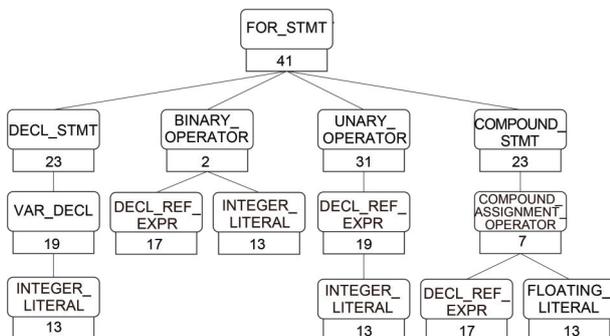


Fig. 2 AST and the hash value of each node for code in Listing 1

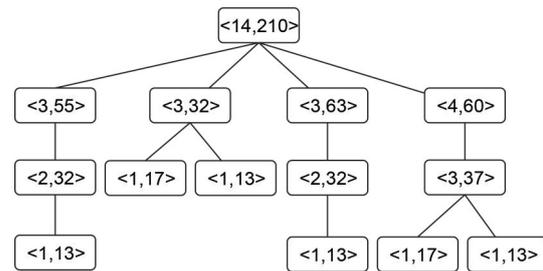


Fig. 3 Hash tree for code in Listing 1, where each node records the sub-tree information rooted at the current node

To speed up the subsequent matching and save memory space, we discard the tree structure and store sub-tree information in the form of lists. The intermediate representation of the code fragment in Listing 1 is shown in Fig. 4. The first dimension of the list represents the number of nodes from 1 to  $N$ , and the hash values of sub-trees with corresponding node numbers are joined into one list. Each hash value list is kept to speed up the search process. We also record the index of code and the total number of nodes in the tree at the start of the list.

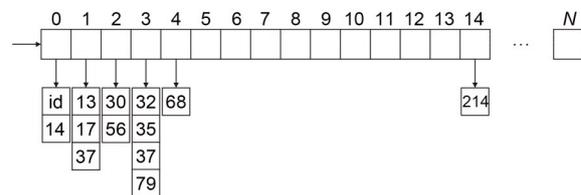


Fig. 4 Intermediate representation of candidate code in Listing 1

### 3.2 Search

In this part, we calculate the similarity scores between the query code and each candidate code, and obtain a candidate code set containing  $K$  code snippets with the highest similarities.

Like candidate code processing, we map the AST of a given query code to a hash tree, similar to Fig. 3. We first filter out the candidate codes that cannot reach the size threshold according to the number of nodes. We traverse the query hash tree in pre-order and search for the hash value of the current node in the candidate hash value list. Since sub-trees with different node numbers represent different structures, only pairs of sub-trees with the same node number need to be compared. For each sub-tree node

in the query hash tree, we obtain the hash value list with the same node number and search for the target hash value using a binary search. If found, the sub-tree represented by the current node is regarded as a match, and all its child nodes will be skipped in the subsequent traversal. We take the total number of nodes under all common sub-trees as the similarity score and record  $K$  code fragments with the highest score. Algorithm 1 shows the process of search, which we will introduce in three parts: preliminary filtering, sub-tree validation, and approximate matching.

### 3.2.1 Preliminary filter based on the node number

When applying code recommendations to real tasks, the vast majority of code snippets in code databases are irrelevant to the query code. However, we need to process all possible pairs of methods to find out potential results. This can be extremely costly, especially on very large datasets. For candidate code fragments whose number of nodes is smaller than that of the query code, their maximum similarity score is the ratio of the number of common nodes to the total number of query nodes. We use size-based heuristics to aggressively eliminate unlikely candidate code snippets upfront. The intuition is that two methods with considerably different sizes are very unlikely to implement the same, or even similar, functionality (Saini et al., 2018). In addition, useful code recommendations are usually larger than query code so that developers can obtain a reference from the extra part. So, for each candidate code, we first judge whether it satisfies Eq. (1):

$$|\text{candidate}_n| > \lambda |\text{query}|, \quad (1)$$

where  $|\text{candidate}_n|$  represents the node number of the current candidate tree,  $|\text{query}|$  represents the node number of the query tree, and  $\lambda$  is an adjustment factor.

If Eq. (1) is satisfied, subsequent matching will be carried out; otherwise, the candidate code will be omitted. This heuristic can lead to some false negatives, especially for code fragments with the same function but differing in texture or structure (Saini et al., 2018). However, in our experiments described below in Section 4.2.3, we observe little or no impact of this on the recall value of the final result under appropriate parameter settings.

---

### Algorithm 1 The search algorithm

---

**Input:** query tree,  $Q$ ; lists of candidate code,  $L$ ; size of target candidate set,  $K$   
**Output:** top- $K$  result list,  $T$

```

1   $V \leftarrow \emptyset$ 
2  for each  $i$  in  $L$ 
3      if sizeMatched( $Q, i$ ) // Filter based on the node number
4           $V \leftarrow V \cup \{<i, \text{traverseTree}(Q, i)>\}$ 
5      end if
6  end for
7   $T \leftarrow$  Rank candidate set  $V$  and obtain the top- $K$  result
  // Traverse the query tree recursively
8  function traverseTree ( $N$ : sub-tree node,  $C$ : list)
9      totalNum  $\leftarrow$  node number of the query tree
10     sub-treeNum  $\leftarrow$  sub-node number of current tree  $N$ 
11     matchedNum  $\leftarrow$  0
12     if isSearched( $N, C[\text{sub-treeNum}]$ )
13         // Validate the sub-tree
14         if validated(totalNum, sub-treeNum,  $N, C$ )
15             return  $N$ 
16         end if
17     end if
18     // If not found, traverse all the sub-trees
19     for each  $j$  in sub-tree of  $N$ 
20         matchedNum  $\leftarrow$  matchedNum + traverseTree( $j, C$ )
21     end for
22     // Approximate matching for the current node
23     if exceedThreshold(matchedNum, sub-treeNum)
24         matchedNum  $\leftarrow$  matchedNum + 1
25     end if
26     return matchedNum
27 end

```

---

### 3.2.2 Sub-tree validation

Since we use only one number to represent the hash value, sub-trees containing different nodes may be mapped to the same hash value, which we call a hash collision. Generally, with the increase of the node number of sub-trees, hash collision is more likely to occur, and when the proportion of sub-trees in the query tree increases, their influence on the final result is more significant. For these sub-trees, additional verification is performed. We adopt a simple but effective method, that is, to verify whether the subsequent nodes of this sub-tree root can also be found in the candidate sub-trees. If they all can be found, the two sub-trees are considered to be matched. We use the following formula to decide the number of nodes to be verified:

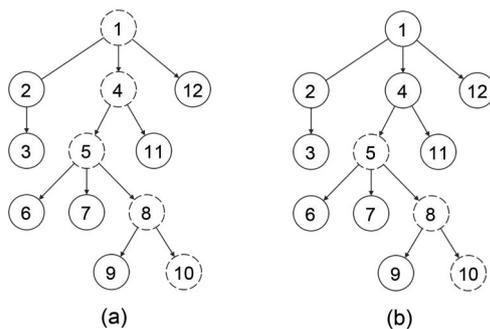
$$T(c, q) = \left\lceil \frac{5c^2}{(c + 10)q} \right\rceil, \quad (2)$$

where  $[t]$  denotes the largest integer no more than  $t$ ,  $c$  denotes the node number of the current sub-tree, and  $q$  denotes the node number of the whole query tree.

We find that almost all false-positive recommending results caused by a hash collision can be eliminated at a low extra cost using this verification algorithm. The corresponding experimental results are given in Appendix B.

### 3.2.3 Approximate sub-tree matching

We can obtain the total number of sub-tree nodes repeated between the query code and each candidate code through the matching process. However, exact sub-tree matching cannot fully reflect the structural similarity. Here, we use the query tree in Fig. 5 as an example. For illustration, we mark the node ID and omit the node number and hash value in this figure.



**Fig. 5** A query tree example in a match without (a) and with (b) approximate sub-tree matching (The solid circle represents a matched node and the dashed circle represents an unmatched node)

Consider a candidate code fragment whose AST is almost the same as that of the query code, the only difference being that node 10 is either missed or replaced by a node with a different hash value. In this instance, all sub-trees containing node 10 obtain different hash values, which means that all sub-tree nodes on the path from node 10 to the root of the query tree cannot be successfully matched (Fig. 5a). Hence, the similarity for this tree is only  $7/12=58.3\%$ , though there is only a different node between the query code and the candidate code. Such minor differences in AST are common in software development (Baxter et al., 1998), and it is reasonable to think that it will be more evident for unfinished code.

So, rather than searching for identical sub-trees, we prefer to find similar ones. We set a similarity

threshold  $T$  to judge whether a sub-tree is matched. For each none-leaf node in the query tree, if no exact match is found in the candidate hash list, additional judgments will be made based on the sub-tree rooted at the current node. That is, if the proportion of matched nodes in the whole sub-tree exceeds  $T$ , the current node will be deemed to be matched. For example, as Fig. 5b shows, if we set  $T$  as 80%, among all nodes on the path from root node 1 to node 10, nodes 1 and 4 are matched while nodes 5, 8, and 10 are unmatched, and the similarity score should be  $9/12=75\%$ , which is a better reflection of code similarity.

### 3.3 Re-ranking

According to our above algorithm, when the number of common leaf nodes is fixed, code fragments with a similar structure to query code can obtain a higher similarity score. In some cases, although some code fragments are quite different from the structure of the query code, they may score higher for similarity than related code because they contain more common sub-trees. Such a situation can easily occur when there is no near-exact matching for the current query in the codebase. In addition, several candidate codes may share the same similarity score, in which case we cannot determine the order of these methods in the recommendation list. To solve these problems and reduce false positive results, we introduce a re-ranking method based on the SW algorithm.

The SW algorithm is used for sequence alignment, which means finding a similar region between two sequences. The purpose of the algorithm is not to compare the whole sequence, but to find similar fragments between the two sequences and determine their similarity.

For query code and each candidate code, we traverse their hash tree in a depth-first order to generate node sequences. Then we use the SW algorithm to obtain the similarity scores between them. The details of SW are presented as follows (Smith and Waterman, 1981):

Assume that  $A=a_1a_2\cdots a_n$  and  $B=b_1b_2\cdots b_m$  are node sequences to be aligned. A similarity is given between sequence elements  $a_i$  and  $b_j$ :

$$s(a_i, b_j) = \begin{cases} \text{match}, & a_i = b_j, \\ \text{mismatch}, & a_i \neq b_j. \end{cases} \quad (3)$$

We set up a matrix  $\mathbf{H}$  whose size is  $(n+1) \times (m+1)$ . Its first row and first column are initialized with 0. Based on the initial value of scoring matrix  $\mathbf{H}$ , the scores of other elements are calculated using the following formula:

$$H_{ij}(i \geq 2, j \geq 2) = \max(H_{i-1, j-1} + s(a_i, b_j), H_{i-1, j} + \text{gap}, H_{i, j-1} + \text{gap}, 0). \quad (4)$$

The highest value of  $H_{ij}$ 's is returned as the similarity score between two sequences.

We combine the SW sequence similarity with the tree-list similarity obtained in the filtering phase as the basis for the re-ranking process, rather than using a single similarity score. This is because the SW algorithm cannot deal with different code statement orders, while tree-list similarity covers such cases and can supplement the final result. The final similarity score is calculated using Eq. (5):

$$\text{Score}(\text{query}, \text{candidate}_n) = \delta \cdot \text{TLsim}(\text{query}, \text{candidate}_n) + \text{SWSim}(\text{query}, \text{candidate}_n), \quad (5)$$

where  $\text{TLsim}$  denotes the tree-list similarity from the search phase,  $\text{SWSim}$  denotes the SW sequence similarity from the re-ranking phase, and  $\delta$  is an adjustment factor.

## 4 Experiments

In this section, we describe our evaluation of the effectiveness of the method proposed in this study. Our experiments were conducted on a 2.60 GHz CPU (Intel i5) PC running Windows 10 OS with 8 GB memory.

### 4.1 Experimental design

#### 4.1.1 Datasets

We used three real-world datasets to evaluate the performance of our method on code recommendation: a Java code repository IJaDataset (Svajlenko and Roy, 2021) and two C code datasets, OJSystem (Mou et al., 2016) and OJNUAA, from pedagogical programming open judge (OJ) systems.

IJaDataset 2.0 is a large Java repository from the SECold Project (<https://sites.google.com/site/aseg-secold/projects/seclone>), containing 25 000 open-source

Java projects. We extracted method-level code fragments and kept those with the number of executable lines between 20 and 30. Then we randomly sampled 300 code fragments and manually created partial code snippets by deleting 30%–50% of code lines from the original code snippets, aiming to check whether our method could produce appropriate recommendations when an exacted match existed in the code database. Three developers with three years of Java project development experience were asked to judge the returned results. If two or more developers thought that the result corresponded to the query code, then we took and recorded it.

OJSystem contains 104 programming problems together with different source codes submitted by students for each problem. There are 500 files submitted under each problem. Two different source codes under programming problems can be regarded as similar if they realize the same functionality. We randomly selected 20 code files as queries under each programming problem and removed them from the candidate set to see if our method could recommend other files under the same problem according to these queries.

We constructed OJNUAA similar to OJSystem. It contains 99 programming problems, each with more than 200 files submitted from a school programming online judgment system. OJNUAA covers a wide range of programming problems from simple to complex, making the average length of code files smaller than that in OJSystem, and the submitted codes for fundamental problems are mostly similar to each other. We also randomly collected 20 code files as queries for each problem on this dataset.

The statistical information of the above datasets (<https://github.com/melond/RecommendationSamples>) is shown in Table 1.

**Table 1 Dataset size information**

Project name	File/Method number	Lines of code
IJaDataset 2.0	718 525	18 237 218
OJSystem	52 000	1 885 262
OJNUAA	34 797	902 133

#### 4.1.2 Metrics

We chose  $\text{Recall}@K$  and the mean reciprocal rank (MRR) as the measurement metrics to evaluate the performance of our proposed method.

Recall@ $K$  is defined as the percentage of query code snippets for which the original method body is found in the top- $K$  methods in the re-ranked search result list. Recall@ $K$  is calculated as

$$\text{Recall@}K = \frac{|\text{Relevance@}K|}{|\text{Queries}|} \times 100\%, \quad (6)$$

where  $|\text{Relevance@}K|$  denotes the number of query code snippets for which the original method body is found in the top- $K$  methods in the result list, and  $|\text{Queries}|$  denotes the number of all queries.

MRR is defined as the average of the reciprocal ranks for all queries, where the reciprocal rank of a query is the inverse of the rank of the first relevant result. A larger MRR value means a higher ranking for the first relevant methods. MRR is calculated as

$$\text{MRR} = \frac{1}{Q} \sum_{q=1}^Q \text{FRank}_q^{-1}, \quad (7)$$

where  $Q$  denotes the total number of query codes, and  $\text{FRank}_q$  denotes the rank of the first relevant result for query  $q$ .

## 4.2 Results

We investigated two research questions (RQs) to validate the effectiveness of our method.

RQ1: How effective is our method in recommending code snippets?

In this study, we attempted to propose a new code recommendation tool. We introduced an approximate sub-tree matching algorithm to search for candidate code sets and used the SW algorithm to re-rank the list. In this RQ, to verify our method's effectiveness, we compared our experimental results with those of two strong baseline works: SENSORY and Aroma.

SENSORY (Ai et al., 2019) is a practical code recommendation tool. It uses the granularity of code statements rather than tokens as input, using the Burrows–Wheeler transform (BWT) algorithm to perform an ordered subsequence search. Since SENSORY is a Java method level tool while OJ databases are constructed by C files, we used Clang to parse the AST, so SENSORY could process the C code.

Aroma (Luan et al., 2018) is a state-of-the-art code recommendation tool. It vectorizes the AST

features and uses a bag-of-words-like strategy to calculate the similarity between the query code and code in the database. It then re-ranks the result list and uses an intersection phase to merge different code snippets in the result list. Since the intersection phase will change the original code snippets from datasets, a manual judgement may be required to assess the correctness of the final result. It is hard to apply to large-scale automated validation and may introduce manual biases. Therefore, we did not adopt the final intersection stage in the experiment, but obtained the experimental results according to the previous part.

Table 2 shows the experimental results of these three methods applied to different datasets. The results show that our method outperformed the compared methods in most situations.

**Table 2 The experimental results from different methods over all datasets**

Dataset	Metric	Value		
		SENSORY	Aroma	Our method
IJaDataset 2.0	Recall@1	100%	99%	99.3%
	Recall@10	100%	100%	100%
	MRR	1	0.995	0.996
OJSystem	Recall@1	67.3%	76.2%	81.1%
	Recall@10	78.6%	86.7%	92.1%
	MRR	0.723	0.821	0.847
OJNUAA	Recall@1	81.5%	91.8%	94.5%
	Recall@10	90.3%	95.5%	97.3%
	MRR	0.909	0.946	0.967

In the experiment on IJaDataset, we treated the results corresponding to the query as a hit, and all three methods achieved good performance. This shows that these methods can make relevant recommendations when there is an exact match in the code dataset in most cases. On OJSystem and OJNUAA, our method performed the best, and SENSORY performed poorly. This is because SENSORY uses a strict code line matching strategy. Its matching process is based on textual information, while the OJ datasets are classified according to functionality. Therefore, it is difficult for SENSORY to recommend the correct result when no exact textually similar candidate code is in the codebase. Aroma chooses to convert the AST into the form of feature vectors and uses feature vectors to calculate the similarity. However, this bag-of-words-like method ignores the

distribution of features in the original AST. Many irrelevant codes will appear with high similarity because they contain similar features to the query code, even though they may not be in the same distribution. Such a phenomenon is more obvious for large code fragments in the codebase and may cause the actual more relevant code to be excluded, thus affecting the final result. In contrast, our method considers the code structure information and stores the candidate codes in the form of sub-tree hash values, so that the candidate codes with a more complete sub-tree structure will have higher similarity. We believe that these are the reasons why our method performed the best on these datasets.

Recommending speed is also important for code recommendation tools. The average time complexity of our method for data processing and matching are  $O(n)$  and  $O(mf)$  respectively, where  $n$  denotes the total AST node number of all candidate codes,  $m$  denotes the node number of the query tree, and  $f$  represents the number of candidate code snippets after size-based filtering. Table 3 shows the recommending time cost of three methods applied to IJaDataset, which contains 718 525 files and 18 237 218 code lines. SENSORY is based on statement sequence matching, so it is time-consuming and not scalable to large datasets. Since Aroma computes similarity based on numerical vectors, it had the best time performance. The time cost of our method was higher than that of Aroma, but still acceptable on such a large dataset.

**Table 3 Time costs of three methods applied to IJaDataset**

Approach	Average recommending time (s)
SENSORY	232.3
Aroma	0.4
Our method	2.1

Our method requires additional storage space to store the hash representation of the AST nodes for code snippets in the code database. Similarly, taking IJaDataset as an example, its source file occupies about 0.63 GB, while the intermediate representation in our method occupies 2.47 GB, which is about four times that of the code source file. The growth of additional space is linearly related to the size of the code database, and it is acceptable for today's data centers.

In summary, our code recommendation method worked well at the method level and file level, and outperformed the baseline methods over all datasets. Also, our method had a good time performance and can be applied to large datasets. Based on the experimental results, we believe that our method can help developers quickly find the code they need.

RQ2: How do the parameters affect the experimental results?

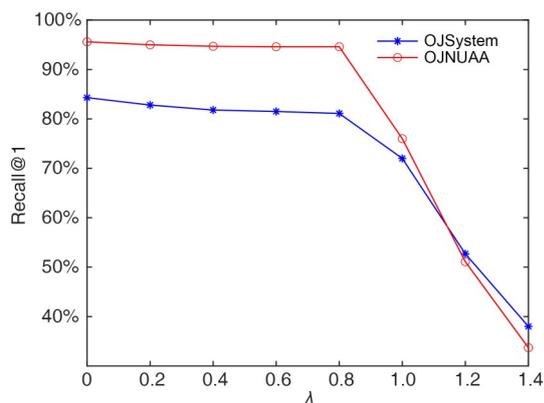
To explore the effects of various parameters on the experimental method, we performed a parameter sensitivity experiment. We tested the two parameters that have the greatest impact on the method:  $\lambda$  in Eq. (1) and  $\delta$  in Eq. (3). The parameters match, mismatch, and gap, required when using the SW algorithm, do not have a great influence on the experimental results (Yang et al., 2018). Therefore, we set them to the empirical values 2, -2, and -1, respectively (Kamalpriya and Singh, 2018).

To answer RQ2, because recommendations should be made repeatedly under each parameter setting in our experiment, and the experiment using IJaDataset needs much manual verification, we chose only OJSystem and OJNUAA for automated validation.

Our method uses Eq. (1) for the size similarity filter. A larger adjustment factor  $\lambda$  can filter out more candidate codes and speed up the overall recommendation, but may also lead to more false negatives. We predefined the range of  $\lambda$  from 0 to 1.4 and used Recall@1 as the metric to estimate the recommending effect.

The experimental results are shown in Fig. 6. The recall rate was close when  $\lambda$  was between 0 and 0.8, and a significant decrease occurred when  $\lambda$  was greater than 0.8. We conclude that the size-based filtering had little or no impact on the final recommendation results when  $\lambda$  was lower than 0.8. At the same time, to reduce unnecessary matching as much as possible and speed up the recommendation, we decided to set  $\lambda$  to 0.8.

Our method uses Eq. (3) to combine tree-list similarity and sequence similarity into one final score. We balanced the influence of two kinds of similarity by parameter  $\delta$ ; e.g., the lower the value of  $\delta$ , the greater the influence of sequence similarity on the final result. We re-ranked the first 50 code segments

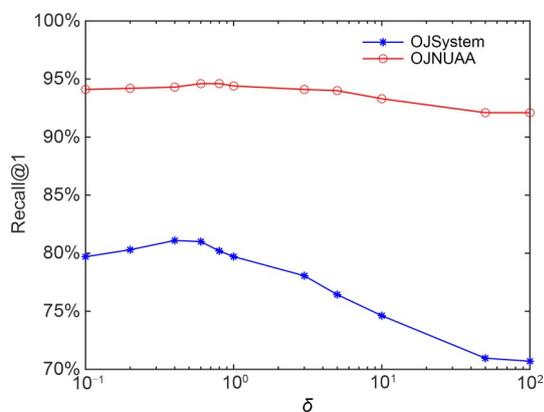


**Fig. 6 Results of our method with different  $\lambda$  values over two datasets**

after the coarse-grained search due to the high time consumption of the SW algorithm. We predefined the range of  $\delta$  from 0.1 to 100, and used Recall@1 as the metric to estimate the recommending effect.

The experimental results over two datasets using Recall@1 are shown in Fig. 7. The recall value using OJNUAA was better than that using OJSystem. This is because OJNUAA covers a broader range and includes many fundamental programming problems. The answers to these basic problems are similar, so the recommendation will be much easier. In addition, the boundaries of some programming problems in OJSystem are not clear, and there are several problems with similar or almost identical functionality, which will also have a specific impact on the results.

Using the OJSystem dataset, our method achieved the best recall results when  $\delta$  was 0.6, and as  $\delta$  continued to grow, the recall rate finally dropped to around 0.7 (Fig. 7). Using OJNUAA, the recall reached the



**Fig. 7 Results of our method with different  $\delta$  values over two datasets**

maximum value when  $\delta$  was set to 0.8. The reason for this difference is that, since the average number of lines of code for each file in OJSystem is larger than that in OJNUAA, sequence similarity plays a more critical role in the final result. The best selection of  $\delta$  was different for different datasets, but when the value of  $\delta$  was between 0.6 and 0.8, our method achieved good results over both datasets. Also, we conclude that the final recommendation results can be significantly improved after the re-ranking phase.

In summary, the best selection of  $\delta$  may change according to the dataset, but our method can achieve good results for both datasets when  $\delta$  is between 0.6 and 0.8.

## 5 Limitations

In this section, we will discuss some limitations associated with our work.

1. The query set. Our method is not being directly implemented by developers in the actual development process because such an experimental setting needs a lot of manual validation and may introduce errors caused by individual biases. Instead, we performed large-scale automated evaluations to test the accuracy of our method. We constructed three real-world code databases and carried out two different types of experiments. On IJaDataset, we tested whether the methods could recommend matching code according to partial code, while for OJSystem and OJNUAA, we evaluated the effectiveness of recommending systems when there may not be a near-exact match in the codebase. Although these two kinds of experimental settings cannot fully cover real development scenarios, they can reflect the performance of the code recommendation system in typical cases.

2. Parameter setting. During the experiment, we set several parameters based on experience or existing experimental results, such as the threshold  $T$  in approximate sub-tree matching and score and penalty parameters in the SW algorithm. We did not conduct complete experiments to analyze the effect of each parameter on the final result, which may have affected the performance of our method to some extent. However, since the experimental results under the current

parameter settings were satisfactory, the parameters in the experimental process should be reasonable. We will fully investigate the effects of each parameter and their combinations in our future work.

## 6 Conclusions and future work

In this study, we have presented a search-based code recommendation technique that combines an approximate sub-tree matching algorithm and a sequence matching algorithm. Specifically, for each candidate code, we first calculate and store the hash value and node number of each sub-tree in their AST. Then, we traverse the AST of the given query code to calculate the similarity between the query code and all the candidate codes in the codebase. Top- $K$  candidate code segments are selected according to a similarity score. We also calculate the AST preorder sequence similarity between the query code and the top- $K$  candidate code using the SW algorithm. Finally, we re-rank the candidate code list based on the two kinds of similarity and return the final recommendation result. To evaluate the effectiveness of our method, we have conducted experiments on several real-world code databases. Experimental results have shown that our method can recommend code with high recall and significantly outperforms compared methods.

In the future, we will consider introducing a clustering process to group functionally similar code fragments together and accelerate the search. We will also consider optimizing the storage structure of the intermediate representation to further reduce the additional storage space required.

### Contributors

Yichao SHAO and Zhiqiu HUANG designed the algorithm, implemented the tool, and drafted the paper. Weiwei LI and Yaosheng YU helped organize the paper. Zhiqiu HUANG and Yaosheng YU revised and finalized the paper.

### Compliance with ethics guidelines

Yichao SHAO, Zhiqiu HUANG, Weiwei LI, and Yaosheng YU declare that they have no conflict of interest.

### References

Ai L, Huang ZQ, Li WW, et al., 2019. SENSORY: leveraging code statement sequence information for code snippets recommendation. Proc IEEE 43<sup>rd</sup> Annual Computer

- Software and Applications Conf, p.27-36.  
<https://doi.org/10.1109/COMPSAC.2019.00014>
- Antunes B, Furtado B, Gomes P, 2014. Context-based search, recommendation and browsing in software development. In: Brézillon P, Gonzalez AJ (Eds.), Context in Computing: a Cross-Disciplinary Approach for Modeling the Real World. Springer, New York, USA, p.45-62.  
[https://doi.org/10.1007/978-1-4939-1887-4\\_4](https://doi.org/10.1007/978-1-4939-1887-4_4)
- Baxter ID, Yahin A, Moura L, et al., 1998. Clone detection using abstract syntax trees. Proc Int Conf on Software Maintenance, p.368-377.  
<https://doi.org/10.1109/ICSM.1998.738528>
- Chen S, Zhang KZ, 2014. An improved algorithm for tree edit distance with applications for RNA secondary structure comparison. *J Comb Optim*, 27(4):778-797.  
<https://doi.org/10.1007/s10878-012-9552-1>
- Đuračik M, Kršák E, Hrkút P, 2020. Searching source code fragments using incremental clustering. *Concurr Comput Pract Exp*, 32(13):e5416. <https://doi.org/10.1002/cpe.5416>
- Holmes R, Murphy GC, 2005. Using structural context to recommend source code examples. Proc 27<sup>th</sup> Int Conf on Software Engineering.  
<https://doi.org/10.1109/ICSE.2005.1553554>
- Jiang H, Nie LM, Sun ZY, et al., 2019. ROSF: leveraging information retrieval and supervised learning for recommending code snippets. *IEEE Trans Serv Comput*, 12(1): 34-46. <https://doi.org/10.1109/TSC.2016.2592909>
- Jiang LX, Mishserghi G, Su ZD, et al., 2007. DECKARD: scalable and accurate tree-based detection of code clones. Proc 29<sup>th</sup> Int Conf on Software Engineering, p.96-105.  
<https://doi.org/10.1109/ICSE.2007.30>
- Kamalpriya CM, Singh P, 2018. Enhancing program dependency graph based clone detection using approximate subgraph matching. Proc IEEE 11<sup>th</sup> Int Workshop on Software Clones, p.1-7.  
<https://doi.org/10.1109/IWSC.2017.7880511>
- Luan SF, Yang D, Barnaby C, et al., 2018. Aroma: code recommendation via structural code search. *Proc ACM Program Lang*, 3:152. <https://doi.org/10.1145/3361527>
- Mou LL, Li G, Zhang L, et al., 2016. Convolutional neural networks over tree structures for programming language processing. Proc 31<sup>st</sup> AAAI Conf on Artificial Intelligence, p.1287-1293.
- Rahman MM, Roy CK, 2014. On the use of context in recommending exception handling code examples. Proc IEEE 14<sup>th</sup> Int Working Conf on Source Code Analysis and Manipulation, p.285-294.  
<https://doi.org/10.1109/SCAM.2014.15>
- Rahman MM, Roy CK, Lo D, 2016. RACK: automatic API recommendation using crowdsourced knowledge. Proc IEEE 23<sup>rd</sup> Int Conf on Software Analysis, Evolution, and Reengineering, p.349-359.  
<https://doi.org/10.1109/SANER.2016.80>
- Sahavechaphan N, Claypool K, 2006. XSnippet: mining for sample code. Proc 21<sup>st</sup> Annual ACM SIGPLAN Conf on Object-Oriented Programming Systems, Languages, and Applications, p.413-430.  
<https://doi.org/10.1145/1167473.1167508>

- Saini V, Farmahinifarahani F, Lu YD, et al., 2018. Oreo: detection of clones in the twilight zone. Proc 26<sup>th</sup> ACM Joint Meeting on European Software Engineering Conf and Symp on the Foundations of Software Engineering, p.354-365. <https://doi.org/10.1145/3236024.3236026>
- Shasha D, Wang JTL, Zhang KZ, et al., 1994. Exact and approximate algorithms for unordered tree matching. *IEEE Trans Syst Man Cybern*, 24(4):668-678. <https://doi.org/10.1109/21.286387>
- Smith TF, Waterman MS, 1981. Identification of common molecular subsequences. *J Mol Biol*, 147(1):195-197. [https://doi.org/10.1016/0022-2836\(81\)90087-5](https://doi.org/10.1016/0022-2836(81)90087-5)
- Svajlenko J, Roy CK, 2021. The mutation and injection framework: evaluating clone detection tools with mutation analysis. *IEEE Trans Soft Eng*, 47(5):1060-1087. <https://doi.org/10.1109/TSE.2019.2912962>
- Yang YM, Ren ZL, Chen X, et al., 2018. Structural function based code clone detection using a new hybrid technique. Proc IEEE 42<sup>nd</sup> Annual Computer Software and Applications Conf, p.286-291. <https://doi.org/10.1109/COMPSAC.2018.00045>
- Ye YW, Fischer G, 2002. Supporting reuse by delivering task-relevant and personalized information. Proc 24<sup>th</sup> Int Conf on Software Engineering. <https://doi.org/10.1109/ICSE.2002.1007995>
- Zhang KZ, Shasha D, 1989. Simple fast algorithms for the editing distance between trees and related problems. *SIAM J Comput*, 18:1245-1262. <https://doi.org/10.1137/0218082>

## Appendix A: Categories of node types for Javalang

**Table A1** Categories of node types for Javalang

Function	Number of node types	Example
Class	10	TypeDeclaration ClassDeclaration InterfaceDeclaration
Type	4	Type BasicType ReferenceType
TypeParameter	1	TypeParameter
Pair and Array	3	Annotation ElementValuePair ElementArrayValue
MemberDeclaration	4	MethodDeclaration FieldDeclaration ConstructorDeclaration
VariableDeclaration	7	ConstantDeclaration VariableDeclaration LocalVariableDeclaration

To be continued

**Table A1**

Function	Number of node types	Example
ExceptionHandlerClause	3	TryResource CatchClause CatchClauseParameter
BasicStatement	7	BlockStatement TryStatement
BranchStatement	2	IfStatement SwitchStatement
LoopStatement	3	DoStatement ForStatement
LoopControlStatement	3	BreakStatement ContinueStatement ReturnStatement
LoopControl	3	SwitchStatementCase ForControl EnhancedForControl

## Appendix B: Recommendation results before and after sub-tree validation

Eq. (2) is used to verify the sub-tree pairs that may produce a hash collision. We tested Recall@1 before and after using hash verification on the OJSys-tem dataset and collected the false-positive results that appeared in the first place in the recommendation list. We analyzed the proportion of recommendation errors caused by hash collision. The experimental results are shown in Table B1.

**Table B1** Recommendation results before and after sub-tree validation

Condition	Recall@1	False positives caused by hash collision
Without validation	75.8%	24.3%
With validation	81.1%	0.3%

Table B1 shows that 24.3% of the false-positive results were caused by hash collisions without the verification phase. After implementing hash validation using Eq. (1), we observed very few cases caused by hash collision. We conclude that the hash verification process can largely eliminate the negative impact of hash collision.