# Enriching Query Semantics for Code Search with Reinforcement Learning

Chaozheng Wang[a], Zhenhao Nong[a], Cuiyun Gao[a,*], Zongjie Li[a], Jichuan Zeng[b], Zhenchang Xing[c] and Yang Liu[d]

[a]*School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, China*

[b]*Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, China*

[c]*Research School of Computer Science, Australian National University, Australia*

[d]*School of Computer Science and Engineering, Nanyang Technology University, Singapore*

## ARTICLE INFO

## ABSTRACT
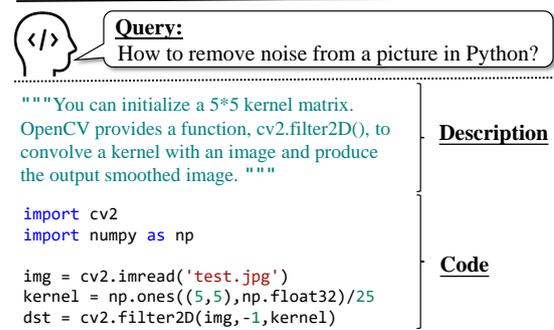
Code search is a common practice for developers during software implementation. The challenges of accurate code search mainly lie in the knowledge gap between source code and natural language (i.e., queries). Due to the limited code-query pairs and large code-description pairs available, the prior studies based on deep learning techniques focus on learning the semantic matching relation between source code and corresponding description texts for the task, and hypothesize that the semantic gap between descriptions and user queries is marginal. In this work, we found that the code search models trained on code-description pairs may not perform well on user queries, which indicates the semantic distance between queries and code descriptions. To mitigate the semantic distance for more effective code search, we propose QueCos, a **Que**ry-enriched **Co**de **s**earch model. QueCos learns to generate *semantic enriched queries* to capture the key semantics of given queries with reinforcement learning (RL). With RL, the code search performance is considered as a reward for producing accurate semantic enriched queries. The enriched queries are finally employed for code search. Experiments on the benchmark datasets show that QueCos can significantly outperform the state-of-the-art code search models.

## 1. Introduction

Searching large corpus of existing source code is a common behavior for developers during software programming. The goal of code search is to retrieve code snippets that most closely match a developer's query, which is generally described in natural language (e.g., the query illustrated in the top of Figure 1). Existing code search approaches can be divided into two categories: information retrieval (IR)-based (e.g., [16, 21, 20]) and deep learning (DL)-based (e.g., [24, 4, 26, 29]). For example, Linstead et al. develop Sourcerer which retrieves similar code snippets based on software textual content and structural information [16]. The IR-based approaches rely on the overlapping tokens or language structures between natural language texts and code snippets, thus suffering from mismatches between the two heterogeneous sources [21]. Recent studies resort to deep learning techniques to remedy the issue by embedding source code and textual code descriptions into the same semantic space. The descriptions are usually written by developers to depict the functions of code snippets (as shown in Figure 1), and can facilitate program comprehension. For example, Zhu et al. propose OCoR to capture the character-



**Figure 1:** An example of a triple of query, description and code.

level and word-level overlaps between code and descriptions based on convolutional network and self attention mechanism [31].

Although promising results were achieved by the deep learning-based approaches, the approaches mostly focused on learning the semantic matching relations between source code and descriptions, and ignored the knowledge gap between input queries and code descriptions. As shown in Figure 1, the code search models trained on pairs of code and descriptions can retrieve the code snippet given the description, since there exist semantically-related keywords between the code and description such as "*cv2*" and "*filter2D*". However, given the user query - "*How to remove noise from a picture in Python?*", which is relatively shorter than the description and limited in context, it will be difficult for the trained models to accurately retrieve the correspond-

---
*Corresponding author

✉ wangchaozheng@stu.hit.edu.cn (C. Wang);
nongzhenhao@stu.hit.edu.cn (Z. Nong);
gaocuiyun@hit.edu.cn (C. Gao);
lizongjie@stu.hit.edu.cn (Z. Li);
jczeng@cse.cuhk.edu.hk (J. Zeng);
zhenchang.xing@anu.edu.au (Z. Xing);
yangliu@ntu.edu.sg (Y. Liu)
ORCID(s):

ing code snippet. To illustrate the semantic gap between queries and descriptions, in this work, we crawled 11,252 Java and 26,237 Python-related code-description-query triples from GitHub and Stack Overflow, one popular question-answering site for developers, respectively. As far as we know, we are the first to prepare such dataset. We found that the lengths of queries are generally shorter than those of the corresponding code descriptions, *e.g.*, 11.97 v.s. 73.61 words on average in Python (as depicted in Table 1 in Section 3). Preliminary experimentation further indicates that employing code-description pairs as training dataset possibly do not generalize well for user queries, which motivates us to enrich queries for effective code search.

There also exists prior work on query expansion to enrich the semantic information of user queries [23, 17, 9]. For example, Liu et al. predict sets of keywords to expand extremely short queries (e.g., queries with two or three words) [17]. The approaches heavily rely on the relevancy of the extended words to the queries and do not explicitly consider the code search performance during expanding the queries. In this paper, we propose to naturally use the search performance as a reward for capturing the query semantics through reinforcement learning (RL). The proposed approach is named as **QueCos**, an abbreviation of **Que**ry-enriched **Co**de **s**earch model. Specifically, QueCos learns to generate the corresponding code descriptions given user queries, and the generated descriptions (called *semantic enriched queries* in the paper) are utilized for improving the code search performance.

In summary, we make the following contributions.

- We prepare the first dataset containing triples of code, description texts and corresponding search queries. Based on the dataset, we observe that the code search models training on code-description pairs may not perform well on user queries.

- We propose QueCos from a novel perspective of generating semantic enriched queries for code search. QueCos naturally captures the key semantic information of queries with the search performance as a reward. The semantic enriched queries are finally utilized to improve the code search task.

- Extensive experiments demonstrate the effectiveness and flexibility of the proposed model on benchmark datasets. The source code and collected datasets are publicly available through this link[1].

**Paper structure.** The remainder of the paper is organized as follows. Section 2 presents our proposed framework. We introduce the experimental datasets, evaluation metrics and baselines in Section 3, and elaborate on the comparison results in Section 4. Section 5 illustrates the related work. We conclude and mention future work in Section 6.

---
[1]Hidden url link

## 2. Methodology

The overview of the proposed approach QueCos is depicted in Figure 2. As can be seen, QueCos mainly includes three components, i.e., code search model, query semantics enrichment model, and hybrid ranking. Initially, a code search model is trained using large code-description pairs (Section 2.1). Then, a query semantic enriching model is designed to generate the corresponding descriptions given user queries based on our collected dataset (Section 2.2), during which RL is adopted to enable the code snippets retrieved by the generated descriptions to be ranked higher. The generated descriptions are treated as semantic enriched queries and not necessary to be exactly close to the descriptions in the ground truth. Finally, both the semantic enriched queries and original queries are employed for the ultimate code search (Section 2.3). We elaborate on the details of each component in the following.

### 2.1. Code Search (CS)

The code search component aims at learning a unified vector representation of both code snippets and descriptions. Since the focus of the work is to mitigate the semantic gap between natural language queries and descriptions for code search, we directly adopt the common framework of the existing code search models [7, 4]. Let $\langle C, D \rangle$ be the set of code-description pairs, where $C$ and $D$ denote the sets of code snippets and descriptions, respectively. The pipeline of existing code search frameworks [7, 10, 31] are shown in the leftmost part of Figure 2, highlighted in green background. Generally, two neural networks are designed separately to encode code and descriptions for mapping the vector representations of the two input sources into the same space, where the neural networks can be Long Short-Term Memory (LSTM), Convolutional Neural Network (CNN) or Transformer [10].

The CS model is trained by minimizing a ranking loss. Specifically, for each description $d \in D$ in the training corpus, the triple of $\langle d, c+, c- \rangle$ is prepared as training instance, where $c+ \in C$ and $c- \in C$ denote the correct code snippet that matches $d$ and sampled negative code snippet that does not match the description (which is randomly sampled from the corpus $C$), respectively. The objective ranking loss function is defined as:

$$\mathcal{L}(\theta)_{CS} = \sum_{<D,C+,C->} \max(0, \epsilon - sim(\mathbf{d}, \mathbf{c}+) + sim(\mathbf{d}, \mathbf{c}-)), \quad (1)$$

where $\theta$ denotes the parameters in the CS model, $\epsilon$ is a constant margin ($\epsilon$ is fixed as 0.05 in all the experiments similar to [7]), $\mathbf{d}$ denotes the description vector of the description $d$, and $\mathbf{c}+$ and $\mathbf{c}-$ denote the code vectors of $c+$ and $c-$, respectively. The *sim* indicates the similarity measurement method which varies for different code search models.

### 2.2. Query Semantic Enrichment (QSE)

The query semantic enrichment (QSE) in QueCos targets at generating semantic-augmented queries given input
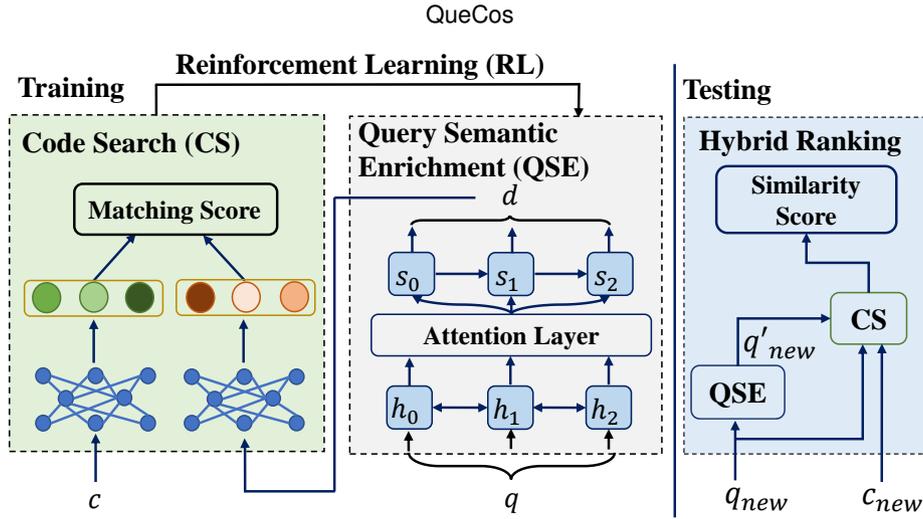
**Figure 2:** Overall framework of the proposed QueCos.

queries. The model is trained to generate corresponding human-provided descriptions based on the queries. Unlike existing sequence-to-sequence models [29], the generated descriptions (i.e., enriched queries) in QueCos are not necessary to be exactly close to the human-provided ones, and the purpose of the training is to capture the main semantics of the queries.

### 2.2.1. QSE model

Given a natural language query $q \in Q$, where $Q$ is the set of collected queries, the QSE learns to produce the corresponding description tokens $d = (d_1, ..., d_N)$.

$$p(\mathbf{d}|\mathbf{q}) = p(d_1|d_0, \mathbf{q}) \prod_{t=2}^{N} p(d_i|d_{1...t-1}, \mathbf{q}) \qquad (2)$$

We adopt standard bi-directional Long Short-Term Memory (LSTM) model with the attention mechanism [1] as QSE model structure to generate descriptions.

$$h_t = \text{Att-Bi-LSTM}(\overrightarrow{h_{t-1}}, \overleftarrow{h_{t+1}}, q_t),$$
$$p(d_t|d_{0...t-1}, \mathbf{q}) = \text{softmax}(\mathbf{W}h_t + b) \qquad (3)$$

where $h_t$ is the decoder hidden state at step $t$, $\mathbf{W}$ and $b$ are learnable weights to project the hidden state $h_t$ to the description vocabulary space.

### 2.2.2. Training QSE via RL

The reinforcement learning (RL) component is designed to render the generated queries in the QSE component able to rank the relevant code snippets ahead while close to the semantics of the descriptions. We view the description generation process as a Markov Decision Process (MDP) [3], and then use *advantage Actor-Critic* or *A2C* algorithm [22] to update the policy gradients. The MDP mainly consists of four ingredients, including state, action, reward and policy.

**State:** In the decoding process, a state maintains the input query $\mathbf{q}$ and the generated tokens $d_{1...t-1}$. We take the

hidden state vector $s_t$ as the vector representation of the state at the $t$-th decoding step.

**Action:** The action for the QSE is to choose the next word $d_t$ from the pre-defined vocabulary. So the action space in our formulation is the vocabulary.

**Reward:** The reward is used to evaluate the quality of the generated descriptions. In this work, the QSE is rewarded based on whether the correct code snippets retrieved by the semantic enriched queries are ranked higher and the relevancy between the semantic enriched queries and the code descriptions in the ground truth. Hence, we define the reward function $r(s_t, d_t)$ at each time step $t$ as:

$$r(s_t, d_t) = \begin{cases} \alpha * \text{Rank}(C, d_{1...t-1}) + \\ (1 - \alpha) * \text{BLEU}(d_{1...t-1}, \mathbf{d}^g), & \text{if } d_t = <\text{EOS}> \\ 0 & \text{otherwise} \end{cases} \qquad (4)$$

where Rank is defined as the popular ranking metric Mean Reciprocal Rank [7] (detailed in Section 3) value, which measures the ranking scores of the correct code snippets given the generated texts. The BLEU [1] score is the common metric for quantitatively estimating the relevance between generated texts $d_{1...t-1}$ and the descriptions $\mathbf{d}^g$ in the ground truth. We use the BLEU-4 score in the paper. The $\alpha$ is the reward tunning parameter and $\alpha \in [0, 1]$.

**Policy:** The policy function computes the probability of choosing $d_i$ as the next token. In this work, we use the policy gradient method [25] to optimize the policy function.

### 2.2.3. Optimization

In order to stabilize the training process, we resort to the A2C algorithm [22]. Specifically, the gradient function is defined as:

$$\nabla \mathcal{L}(\phi) = \mathbb{E}\left[\sum_{t=1}^{N} (R_t(s_t, d_t) - V(s_t)) \nabla \log P(d_t|d_{1...t-1}, C; \phi)\right],$$

**Table 1**

Statistics of the CodeSearchNet dataset and our collected dataset. Our collected dataset contains aligned code, description texts, and queries. The statistics include the number of pairs and the average lengths of code snippets, descriptions and queries. "-" means there are no such data in the dataset.

| Dataset | Training Set | Validation Set | Test Set | Avg. Code Token Len. | Avg. Desc. Token Len. | Avg. Query Token Len. |
|---------|-----------|-----------|---------|--------------------|--------------------|--------------------|
| Java (CodeSearchNet) | 450,941 | 15,053 | 26,717 | 162.41 | 38.87 | - |
| Java (Ours) | 9,015 | 1,117 | 1,120 | 506.34 | 55.75 | 12.91 |
| Python (CodeSearchNet) | 409,230 | 22,906 | 22,104 | 167.51 | 54.21 | - |
| Python (Ours) | 21,009 | 2,596 | 2,632 | 131.21 | 73.61 | 11.97 |

$$\tag{5}$$

where $R_t(s_t, d_t) = \sum_{t' \geq t} r(s_{t'}, d_{t'})$ is the return for generating word $d_t$ given state $s_t$, $V(s_t)$ is the state value function that estimates the future reward given $s_t$, we construct a critic network to predict $V(s_t)$ and optimize it with a mean square error (MSE) loss $\mathcal{L}(\rho) = \mathbb{E}_{D \sim P(\cdot|Q)}[\sum_{t=1}^{|D|}(V(s_t; \rho) - R(Q, D))^2]$.

## 2.3. Hybrid Ranking

In the training phase, we first train the code search component and then utilize the RL component to generate semantic enriched queries for the code search task. In the testing phase, we design a hybrid ranking component to consider both the original queries and generated enriched queries for returning the ultimate search results.

Specifically, for each input query $q$, we derive the semantic enriched queries $q'$ via the QSE component. The final similarity matching score between the query $q$ and candidate code snippet $c$ is computed as:

$$score(\mathbf{q}, \mathbf{c}) = \beta * sim(\mathbf{q}', \mathbf{c}) + (1 - \beta) * sim(\mathbf{q}, \mathbf{c}), \tag{6}$$

where $\beta$ is a parameter to be adjusted experimentally, $0 \leq \beta \leq 1$. $\mathbf{q}$, and $\mathbf{q}'$, and $\mathbf{c}$ denote the encoded vectors of $q$, $q'$, and $c$, respectively.

## 3. Experimental Setup

### 3.1. Dataset

#### 3.1.1. CodeSearchNet.

CodeSearchNet [10] is a publicly-available GitHub repository. It provides thousands of pairs of code snippets and the corresponding natural language descriptions. We conducted preprocessing following Gu's work [7] by first splitting source code tokens of the form CamelCase and snake case to respective sub-tokens, and then taking lowercase. We removed empty items afterwards. The numbers of the subject code-description pairs for the Java and Python datasets are shown in Table 1.

#### 3.1.2. Our Collected Datasets.

Since no dataset containing aligned code snippets, descriptions, and queries were released, we decided to prepare such dataset. The dataset was crawled based on the SOTorrent dataset [2]. SOTorrent comprises pairs of SO (Stack Overflow) posts and GitHub files, where the GitHub files have referred to the corresponding posts, i.e., the code descriptions in the GitHub files clearly cite the post URLs[2]. Since we focus on the code written in Java and Python, we filtered the code snippets in other languages out. In SOTorrent, only links of the SO posts and GitHub files were provided, so we designed a crawler to traverse all the GitHub links, locate and save the code snippets referred to the SO posts and the corresponding code descriptions. To capture the SO queries, we accessed the SO posts via the SO links and obtained the queries according to the HTML element. The preprocessing procedure for the crawled data is also similar to [7]'s work. We finally split the preprocessed data to be training set, valid set and test set by 8:1:1. The statistics of the subject data are illustrated in Table 1.

### 3.2. Evaluation Metric

Following [7, 10, 31], we evaluate the model performance with the standard metrics, including R@k and MRR (Mean Reciprocal Rank).

#### 3.2.1. R@k

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

$$R@k = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \delta(FRank_q \leq k), \tag{7}$$

where $Q$ denotes the query set and $FRank_q$ denotes the rank of the correct answer for query $q$. The function $\delta(Frank_q \leq$

---

[2]Developers write the SO post URLs in the code descriptions for future program comprehension.

$k$) ∈ {0, 1} returns 1 if the rank of the correct answer within the top $k$ returning results otherwise returns 0. A higher $R@k$ indicates a better code search performance.

### 3.2.2. MRR

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

$$M\,RR = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \frac{1}{F\,Rank_q}. \tag{8}$$

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

### 3.3. Baselines

Since the code search component in QueCos can be built upon any existing code search models, we mainly choose two popular approaches and one state-of-the-art approach for evaluation. **DeepCS** [7] is one popular code search model, which employs neural networks such as RNNs to embed both code snippets and descriptions into a joint vector space. Besides the code tokens, it also considers API sequences and method names for learning the representations of code snippets. Since the API sequences are hard to be extracted for the experimental datasets, we slightly modify the original model to only combine method names and code tokens for learning the code representations. **UNIF** [4] is similar to DeepCS but with attention mechanism involved. Besides, UNIF only considers code tokens during embedding code snippets. **OCoR** [31] is one of the state-of-the-art models. It captures the overlaps between code and descriptions in two levels, including character level and word/identifier level, for calculating the similarities between code and descriptions.

We also choose one popular query expansion approach, denoted as **QE** [19], as baseline. QE utilizes WordNet [13] to extend the queries with synonyms for effective code search. It can also be built upon any code search models.

### 3.4. Implementation Details

Following the prior studies [7, 10, 31], we tokenized code snippets, descriptions, and queries. Specifically, code snippets are split in camel case and snake case. Then we involved the top 10,000 words in the training set as the vocabularies of code snippets, descriptions, and queries, respectively, according to the word frequencies. In addition, we used two symbols to represents the beginning and end of a sentence for sentence generation in QSE. The word embeddings in both CS and QSE models are randomly initialized. The numbers of hidden states in the encoder and decoder are both defined as 256. The parameter $\alpha$ (in Equ. 4) to tune the reward function is set as 1.0, and the parameter $\beta$ (in Equ. 6) for adjusting the hybrid ranking score is defined as 0.6. Detailed analysis about the parameter settings will

be introduced in Section 4.3. We use the Adam optimizer with a learning rate 1e-3 and the learning rate decay ratio is defined as 0.5. For the baseline models, we use the same model hyper-parameters as the original papers. For evaluation, we fix a set of 999 negative code snippets $c$- for each test pair $\langle q, c+ \rangle$ following [10].

During the training phase, we first train the CS model for 120 epochs, and then train the QSE model for another 20 epochs. Subsequently, we start to train the reinforcement learning component by pretraining the critic network for 10 epochs and jointly training the critic and actor networks for 40 epochs. All the experiments run on a server with 8 * Nvidia Tesla P100 and each one has 16GB graphic memory. Depending on the involved CS model in QueCos, the training process lasts from 70 to 240 hours.

## 4. Results

Based on the preceding experimental setup, we first illustrate the semantic gap between queries and descriptions, and then evaluate the effectiveness of the proposed QueCos model from three aspects, including comparison with the baselines, parameter analysis, and case study.

### 4.1. Preliminary Experiments

Here we investigate whether there exists semantic gap between user queries and code descriptions. Specifically, we analyzed the performance of the model trained with code-description pairs on different types of evaluation sets, including queries and descriptions. We choose the three baseline models, i.e., DeepCS, UNIF, and OCoR for experimentation. All the models are trained on the code-description pairs provided by CodeSearchNet [10], and then evaluated with the crawled query test set and description test set in Table 1. The results are illustrated in the Table 2. We can observe that using queries as test set presents lower performance than using descriptions as test set for all the three code search baselines. For example, the MRR values decrease by 12.8% and 32.9% on average on the Java and Python datasets, respectively; and the R@1 values also show a downtrend, with the average decreasing rates at 21.8% and 42.9% on the two datasets, respectively. The results imply the semantic gap between queries and description texts. Thus, models trained on code-description pairs possibly do not perform well for practical user queries, and mitigating the semantic gap between queries and descriptions is necessary for more accurate code search.

### 4.2. Main Comparison Results

We compare QueCos with the four baseline models introduced in Section 3 on the collected datasets. Since the code search component in QueCos and QE can be built upon any code search models, we analyze the performance with each of the code search baselines involved. The comparison results are illustrated in Table 3. The approach "QueCos w/o RL" indicates the proposed QueCos model without the RL procedure, that is, the code search performance is

**Table 2**
Performance of the baseline models using different evaluation sets.

| Approach | Evaluation Set | Java | | | | Python | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | R@1 | R@5 | R@10 | MRR | R@1 | R@5 | R@10 | MRR |
| DeepCS | Description | 0.202 | 0.389 | 0.488 | 0.297 | 0.160 | 0.298 | 0.368 | 0.231 |
| | Query | 0.180 | 0.381 | 0.483 | 0.278 | 0.100 | 0.223 | 0.287 | 0.161 |
| UNIF | Description | 0.213 | 0.450 | 0.561 | 0.324 | 0.130 | 0.276 | 0.347 | 0.202 |
| | Query | 0.156 | 0.439 | 0.568 | 0.285 | 0.078 | 0.202 | 0.261 | 0.143 |
| OCoR | Description | 0.292 | 0.561 | 0.654 | 0.412 | 0.230 | 0.423 | 0.507 | 0.322 |
| | Query | 0.211 | 0.460 | 0.583 | 0.330 | 0.112 | 0.282 | 0.363 | 0.196 |

**Table 3**
Comparison results with baseline models. The bold figures indicate the best results.

| Approach | Java | | | | Python | | | |
|---|---|---|---|---|---|---|---|---|
| | R@1 | R@5 | R@10 | MRR | R@1 | R@5 | R@10 | MRR |
| **DeepCS-based** | | | | | | | | |
| DeepCS | 0.180 | 0.381 | 0.483 | 0.278 | 0.100 | 0.223 | 0.287 | 0.161 |
| QE | 0.169 | 0.361 | 0.464 | 0.266 | 0.105 | 0.222 | 0.288 | 0.167 |
| QueCos w/o RL | 0.110 | 0.239 | 0.322 | 0.180 | 0.048 | 0.131 | 0.180 | 0.095 |
| QueCos w/o HR | 0.192 | 0.402 | 0.498 | 0.294 | 0.117 | 0.252 | 0.309 | 0.182 |
| QueCos (ours) | **0.282** | **0.519** | **0.611** | **0.394** | **0.149** | **0.297** | **0.370** | **0.224** |
| **UNIF-based** | | | | | | | | |
| UNIF | 0.156 | 0.439 | 0.568 | 0.285 | 0.078 | 0.202 | 0.261 | 0.143 |
| QE | 0.155 | 0.431 | 0.553 | 0.281 | 0.072 | 0.179 | 0.250 | 0.133 |
| QueCos w/o RL | 0.098 | 0.296 | 0.408 | 0.197 | 0.066 | 0.168 | 0.231 | 0.122 |
| QueCos w/o HR | 0.159 | 0.450 | 0.565 | 0.293 | **0.128** | 0.293 | 0.359 | 0.207 |
| QueCos (ours) | **0.199** | **0.517** | **0.660** | **0.347** | 0.125 | **0.319** | **0.394** | **0.217** |
| **OCoR-based** | | | | | | | | |
| OCoR | 0.211 | 0.460 | 0.583 | 0.330 | 0.112 | 0.282 | 0.363 | 0.196 |
| QE | 0.227 | 0.511 | 0.647 | 0.358 | 0.115 | 0.291 | 0.380 | 0.202 |
| QueCos w/o RL | 0.183 | 0.418 | 0.504 | 0.289 | 0.095 | 0.227 | 0.305 | 0.163 |
| QueCos w/o HR | 0.216 | 0.516 | 0.613 | 0.349 | **0.184** | **0.408** | **0.492** | 0.282 |
| QueCos (ours) | **0.242** | **0.553** | **0.649** | **0.375** | **0.184** | 0.398 | 0.489 | **0.289** |

not considered as reward for the query enrichment. The approach "QueCos w/o HR" denotes the QueCos without the HR component, i.e., the semantic enriched queries generated in QSE component are directly used for code search. Based on the comparison results, we achieve the following observations:

**Observation 1: QueCos significantly improves the baseline models given user queries.** As can be seen in Table 3, QueCos presents better performance than the corresponding code search model on both Java and Python datasets. For example, QueCos increases the accuracy of the corresponding code search model by 33.0%, 24.7%, 23.8% and 21.9% in terms of R@1, R@5, R@10 and MRR, respectively, for the Java dataset on average. The results demonstrate that with the query semantics augmented, QueCos can more accurately recommend relevant code snippets for given user queries. It is unsurprising that QueCos with OCoR achieves the best results on the Python dataset since OCoR

already presents the best performance among the three baselines. But we also observe that QueCos with DeepCS shows the best results on the Java dataset in terms of R@1 and MRR, which further highlights the effectiveness of QueCos for enriching user queries for the task.

**Observation 2: The query semantic enriching component is more effective than the popular query expansion approach.** By comparing QueCos with the popular query expansion approach QE, we can discover that QueCos can significantly improve the performance of QE. However, with QE involved, the code search models may perform worse, for example, QE with DeepCS degrades the performance of DeepCS by 6.11% in terms of R@1 on the Java dataset. This indicates that explicitly complementing the queries with synonyms may not contribute the model performance. The results also prove the usefulness of the query semantic enriching component in QueCos for code search.

**Observation 3: The RL and HR components in Que-Cos plays a key role in enriching the query semantics.** Comparing QueCos w/o RL with QueCos, we can find that QueCos performs worse without the RL component, presenting even poorer performance than the corresponding base code search model. The phenomenon indicates that simply using the generated descriptions may bias the semantics of the queries and lead to inaccurate code search. This is reasonable since the queries are generally shorter than the descriptions according to Table 1, and the model is hard to accurately capture the semantic relations between queries and descriptions through merely training on query-description pairs. We provide an error case in Section 4.5 for further illustration.

Comparing QueCos w/o HR with QueCos, we can observe that the effectiveness of only using the generated enriched queries for code search is limited. Integrating the original queries could further enhance the model performance.

### 4.3. Parameter Analysis

In the section, we analyze the impact of two important hyper-parameters, including the reward tuning parameter $\alpha$ and the hybrid ranking parameter $\beta$, on the performance of QueCos. Figure 3 and Figure 4 illustrate the performance variations of QueCos with different base code search models as the parameters changes, respectively. During analyzing the impact of $\alpha$, we did not conduct the parameter analysis for the OCoR-based QueCos due to the huge computation resources required by OCoR[3]. As can be seen in Figure 3, when $\alpha$ increases, the model performance presents an obvious uptrend in terms of both R@1 and MRR. The results indicate that considering the similarity between the generated texts and the descriptions in the ground truth as a reward during semantic enriched query generation is unhelpful for the code search task. This is reasonable since only involving the code ranking performance as a reward can render the RL component focus on generating the descriptions that could rank relevant code snippets higher. Thus, we define the reward tuning parameter $\alpha$ as 1.0 during the experiments, i.e., the relevancy between the generated text and the human-provided description is not considered as a reward.

According to Figure 4, The performance variations of DeepCS-based QueCos along with the value are close to an inverted "U" shape, but the performance variations of UNIF-based or OCoR-based QueCos present an increasing trend along with the value. From Table 1, we observe that the code descriptions in the Python dataset are relatively longer than those in the Java dataset. We guess the generated enriched queries could deliver more accurate semantics for the Python dataset, so the relevance between the enriched queries and code snippets is more important for the search results comparing to the Java dataset. We fix the value of $\beta$ as 0.6 in the experiments since the model would approach the optimal accuracy.

---

[3]The training process of OCoR lasts around 240 hours.

| Query | Java compare two lists' object values? |
|---|---|
| Description in Ground Truth | NOTE: will not work with duplicate objects |
| Query Enriched by QE | comparison compare equivalence liken equate two 2 ll deuce ii list listing tilt inclination lean name heel number object... |
| Query Enriched by QueCos | converts list length of code taken its returned list code list> sorts a list using widths as a list integers |

**Table 4**
User query, description in the ground truth, and the enhanced query by QE and QueCos for the successful case 1.

### 4.4. Case Study

```
1  private Boolean equalLists(List<TariffSpecification>
   ↪ listA, List<TariffSpecification> listB) {
2      return listA.size() == listB.size() &&
       ↪ listA.containsAll(listB);
3  }
```

Code Listing 1: Successful case 1, with more details listed in Table 4.

```
1  def _add_subelem(root_element, name, value):
2      if value is None:
3          return
4      if type(value) is dict:
5          if name == "link":
6              ET.SubElement(root_element, name, value)
7          elif name == 'content':
8              e = ET.Element(name, type =
               ↪ value['type'])
9              e.append(CDATA(value['content']))
10             root_element.append(e)
11         else:
12             subElem = ET.SubElement(root_element,
               ↪ name)
13             for key in value:
14                 _add_subelem(subElem, key,
                   ↪ value[key])
15     else:
16         ET.SubElement(root_element, name).text =
           ↪ value
```
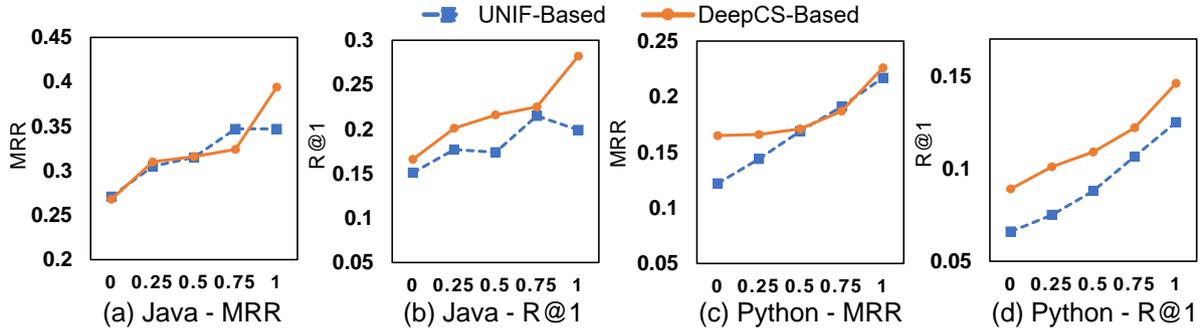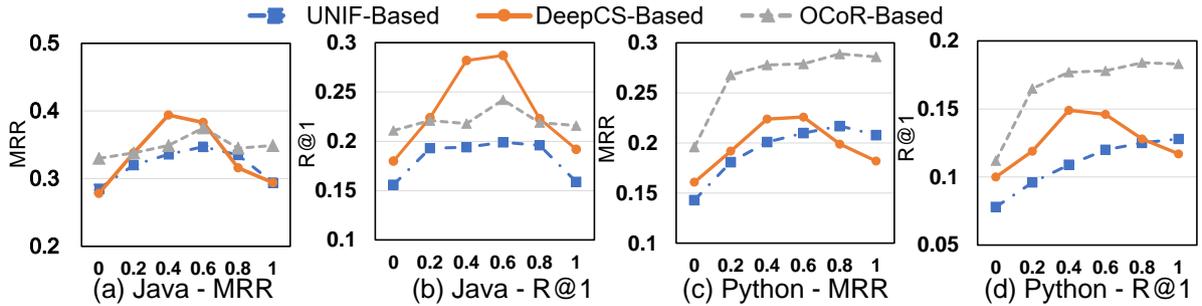
Code Listing 2: Successful case 2, with more details shown in Table 5.

Code Listing 1 and Table 4 show a case in which Que-Cos with UNIF-based successfully returns the correct code snippet of the query "Java compare two lists' object values?" at the top of the ranking list, and the basic UNIF model only ranks the code snippet at 11. We can observe that the code description does not indeed reflect the functionality of the code, i.e., comparison of two lists' objective values. The generated semantic enriched query of Que-Cos is "converts list length of code taken its returned list code list> sorts a list using widths as a list integers", which could deliver more semantics than the input query although the grammar is not exactly correct. The expanded query by the QE approach is "comparison compare equivalence liken equate two 2 ll deuce ii list listing...", which also only ranks

QueCos



**Figure 3:** Parameter sensitive study on the reward tuning parameter $\alpha$. The parameter analysis for OCoR-based QueCos was ignored due to the huge computation resources of OCoR.



**Figure 4:** Parameter sensitive study on the hybrid ranking parameter $\beta$.

| | |
|---|---|
| Query | How to output Cdata using element tree? |
| Description in Ground Truth | HORRIBLE HACK! A wee hack too, the content node must be converted to a CDATA block. This is a sort of cheat, see: |
| Query Enriched by QE | end_product output yield output_signal production outturn turnout exploitation victimization victimisation use utilize utilis... |
| Query Enriched by QueCos | xml xml element note retaining elements are quite source extended xml element xml etree node element as value save a tree object from xml to xml map back xml xml |

**Table 5**
User query, description in the ground truth, and the enhanced query by QE and QueCos for the successful case 2.

the code snippet at 11. In addition, QueCos with DeepCS-based and OCoR-based can rank the correct code snippet at 9 and 19, respectively, while the basic DeepCS and OCoR provide the rank at 8 and 30 respectively. Overall, the generated semantic enriched query can be helpful for more accurate code search.

For the example shown in the Code Listing 2 and Table 5, the QueCos with UNIF-based improves the rank of

the correct code snippet from 60 to 2 compared to the basic UNIF approach [4]. The proposed model can capture that the terms "Cdata" and "element tree" in the query are relevant to the XML language, and the generated semantic enriched query clearly includes terms such as "xml", "element", and "etree"; while the extended query by the QE approach fails to catch the relation of the query to XML, returning the code snippet at 16. So we suppose that through enriching the query semantics, QueCos can more accurately rank the code snippets.

### 4.5. Error Study

```
1  SuppressWarnings(""unchecked"")
2  Map<String, String> myMap = (Map<String, String>)
   ↪  deserializeMap();
```

Code Listing 3: An error Case, with more details in Table 6.

Code Listing 3 and Table 6 show an error case. UNIF ranks corresponding code snippet at 4 according to the original query, however, UNIF-based QueCos only ranks it as 10. The poor performance can be attributed to that QueCos misunderstands the meaning of the word "*address*" in the query, and hence the enriched query contains irrelevant keywords, e.g., "*host*", "*url*", "*ip*", etc. In the future, we will design a quality assurance filtering component to distinguish the semantically-irreverent keywords in the generated enriched queries.

| Query | How do I address unchecked cast warnings? |
|---|---|
| Description in Ground Truth | None |
| Query Enriched by QE | address speech destination savoir-faire turn_to speak direct call cover accost unbridled cast mold mould form plaster_cast casting roll... |
| Query Enriched by QueCos | verify the device class wifi on host ( s , ip listener , url . , it will dynamically handle web fixed allowing the passing and run of the context ( such , domain, dynamic ip address. |

**Table 6**
User query, description in the ground truth, and the enhanced query by QE and QueCos for the error case in Code Listing 3.

## 5. Related Work

### 5.1. Code Search

As the rapid growing of open-source code corpus, to retrieve the code fragments satisfying a user's intent with high accuracy concerning developing productivity. Prior studies on code search focus on figuring out the latent relationship between NL queries and code snippets. The work [21] proposes Portfolio, modelling the navigation behavior of programmers with random surfer and returning a chain of functions. Lv et al. [20] propose CodeHow, a code search tool that first retrieves all possible API calls to augment the queries.

Recently, an increasing number of works using neural networks for code search have been proposed [17, 29, 4, 31]. Sachdev et al. present a neural network-based model, which creates continuous vector representations of codes for comparison [24]. In the work [17], Liu et al. present NQE, which predicts the keywords in the NL queries and expands them in a productive way to improve performance for shorter queries. Gu et al. propose DeepCS which embeds both the code snippets and queries into a joint vector space to measure the similarity between them [7]. [4] introduces UNIF, a bag-of-words-based network which converts code snippets and docstring tokens into embedding matrices with supervised learning method. The authors in [29] treat code annotation and code search as two divided tasks and use the generated code annotations to improve the performance of code search. However, no previous studies have explicitly involved the search performance for query semantic enrichment for improving the task. Zhu et al. [31] propose an overlap-aware OCoR, which explicitly considers the overlap between code and descriptions in both word and character levels. Based on the word-level and character-level representations of the code and descriptions, OCoR utilizes MLP to compute the relevance.

### 5.2. Code Representation Learning

Code representation learning aims at learning the semantics of programs for facilitating various downstream tasks related to program comprehension, such as code clone detection, code summarization, bug detection [30, 7, 10, 28, 14, 27, 15, 11], etc. The development of deep learning techniques boosts the research on code representation learning. In the work [7, 10], code snippets are split into tokens and fed into neural networks such as RNNs and multi-head attentions for the representation learning. Considering the structural nature of code, [28, 14, 30] combine the abstract syntax trees (ASTs) into neural networks for capturing the code semantics. LeClair et al. [14] use GNN-based encoder to model the AST of each program subroutine.

Recent work adapts pre-training techniques in the natural language processing field [5, 12, 18] for better program comprehension. For example, the work [6] presents Code-BERT, the first bimodal pre-trained model for programming language (PL) and natural language (NL). Guo et al. [8] later propose GraphCodeBERT, which combines both code tokens and the data flow information during pre-training, and achieves more promising performance.

## 6. Conclusions

In this paper, we have proposed a novel RL-based query-enriched code search model, named QueCos, for more effective code search. QueCos can capture the main semantics of user queries through learning to generate the corresponding descriptions, where RL is adopted to render the generated descriptions to be able to rank the relevant code snippets higher. We also crawled the first dataset containing triples of code, descriptions, and queries. Experiments demonstrate that QueCos is more effective than the baseline models and also flexible to incorporate any code search models. In the future, we will combine external knowledge graph and adapt Transformer to further enrich the semantics of the queries.

## References

[1] Bahdanau, D., Cho, K., Bengio, Y., 2015. Neural machine translation by jointly learning to align and translate, in: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.

[2] Baltes, S., Dumani, L., Treude, C., Diehl, S., 2018. Sotorrent: reconstructing and analyzing the evolution of stack overflow posts, in: Proceedings of the 15th International Conference on Mining Software Repositories, MSR 2018, Gothenburg, Sweden, May 28-29, 2018, pp. 319–330.

[3] Bellman, R., 1957. A markovian decision process. Journal of mathematics and mechanics , 679–684.

[4] Cambronero, J., Li, H., Kim, S., Sen, K., Chandra, S., 2019. When deep learning met code search, in: Dumas, M., Pfahl, D., Apel, S., Russo, A. (Eds.), Proceedings of the ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/SIGSOFT FSE 2019, Tallinn, Estonia, August 26-30, 2019, ACM. pp. 964–974.

[5] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2018. Bert: Pretraining of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 .

[6] Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., Zhou, M., 2020. Codebert: A pre-trained model for programming and natural languages. CoRR abs/2002.08155.

[7] Gu, X., Zhang, H., Kim, S., 2018. Deep code search, in: Proceedings of the 40th International Conference on Software Engineering, ICSE 2018, Gothenburg, Sweden, May 27 - June 03, 2018, pp. 933–944.

[8] Guo, D., Ren, S., Lu, S., Feng, Z., Tang, D., Liu, S., Zhou, L., Duan, N., Svyatkovskiy, A., Fu, S., Tufano, M., Deng, S.K., Clement, C.B., Drain, D., Sundaresan, N., Yin, J., Jiang, D., Zhou, M., 2020. Graph-codebert: Pre-training code representations with data flow. CoRR abs/2009.08366.

[9] Huang, Q., Yang, Y., Cheng, M., 2019. Deep learning the semantics of change sequences for query expansion. Softw., Pract. Exper. 49, 1600–1617. URL: https://doi.org/10.1002/spe.2736, doi:10.1002/spe.2736.

[10] Husain, H., Wu, H., Gazit, T., Allamanis, M., Brockschmidt, M., 2019. Codesearchnet challenge: Evaluating the state of semantic code search. CoRR abs/1909.09436.

[11] Lam, A.N., Nguyen, A.T., Nguyen, H.A., Nguyen, T.N., 2017. Bug localization with combination of deep learning and information retrieval, in: 2017 IEEE/ACM 25th International Conference on Program Comprehension (ICPC), IEEE. pp. 218–229.

[12] Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R., 2019. Albert: A lite bert for self-supervised learning of language representations. arXiv preprint arXiv:1909.11942 .

[13] Leacock, C., Chodorow, M., 1998. Combining local context and wordnet similarity for word sense identification. WordNet: An electronic lexical database 49, 265–283.

[14] LeClair, A., Haque, S., Wu, L., McMillan, C., 2020. Improved code summarization via a graph neural network, in: Proceedings of the 28th International Conference on Program Comprehension, pp. 184–195.

[15] Li, Y., Wang, S., Nguyen, T.N., Van Nguyen, S., 2019. Improving bug detection via context-based code representation learning and attention-based neural networks. Proceedings of the ACM on Programming Languages 3, 1–30.

[16] Linstead, E., Bajracharya, S.K., Ngo, T.C., Rigor, P., Lopes, C.V., Baldi, P., 2009. Sourcerer: mining and searching internet-scale software repositories. Data Min. Knowl. Discov. 18, 300–336.

[17] Liu, J., Kim, S., Murali, V., Chaudhuri, S., Chandra, S., 2019a. Neural query expansion for code search, in: Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages, pp. 29–37.

[18] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V., 2019b. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 .

[19] Lu, M., Sun, X., Wang, S., Lo, D., Duan, Y., 2015. Query expansion via wordnet for effective code search, in: 22nd IEEE International Conference on Software Analysis, Evolution, and Reengineering, SANER 2015, Montreal, QC, Canada, March 2-6, 2015, pp. 545–549.

[20] Lv, F., Zhang, H., Lou, J., Wang, S., Zhang, D., Zhao, J., 2015. Code-how: Effective code search based on API understanding and extended boolean model (E), in: 30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015, Lincoln, NE, USA, November 9-13, 2015, pp. 260–270.

[21] McMillan, C., Grechanik, M., Poshyvanyk, D., Xie, Q., Fu, C., 2011. Portfolio: finding relevant functions and their usage, in: Proceedings of the 33rd International Conference on Software Engineering, ICSE 2011, Waikiki, Honolulu , HI, USA, May 21-28, 2011, pp. 111–120.

[22] Mnih, V., Badia, A.P., Mirza, M., Graves, A., Lillicrap, T.P., Harley, T., Silver, D., Kavukcuoglu, K., 2016. Asynchronous methods for deep reinforcement learning, in: Balcan, M., Weinberger, K.Q. (Eds.), Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016, JMLR.org. pp. 1928–1937.

[23] Nie, L., Jiang, H., Ren, Z., Sun, Z., Li, X., 2016. Query expansion based on crowd knowledge for code search. IEEE Trans. Services Computing 9, 771–783.

[24] Sachdev, S., Li, H., Luan, S., Kim, S., Sen, K., Chandra, S., 2018. Retrieval on source code: a neural code search, in: Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages, MAPL@PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018, pp. 31–41.

[25] Sutton, R.S., Barto, A.G., et al., 1998. Introduction to reinforcement learning. volume 135. MIT press Cambridge.

[26] Wan, Y., Shu, J., Sui, Y., Xu, G., Zhao, Z., Wu, J., Yu, P.S., 2019. Multi-modal attention network learning for semantic source code retrieval, in: 34th IEEE/ACM International Conference on Automated Software Engineering, ASE 2019, San Diego, CA, USA, November 11-15, 2019, pp. 13–25.

[27] Wan, Y., Zhao, Z., Yang, M., Xu, G., Ying, H., Wu, J., Yu, P.S., 2018. Improving automatic source code summarization via deep reinforcement learning, in: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, pp. 397–407.

[28] Wang, W., Li, G., Shen, S., Xia, X., Jin, Z., 2020. Modular tree network for source code representation learning. ACM Transactions on Software Engineering and Methodology (TOSEM) 29, 1–23.

[29] Yao, Z., Peddamail, J.R., Sun, H., 2019. Coacor: Code annotation for code retrieval with reinforcement learning, in: The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019, pp. 2203–2214.

[30] Zhang, J., Wang, X., Zhang, H., Sun, H., Wang, K., Liu, X., 2019. A novel neural source code representation based on abstract syntax tree, in: 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE), IEEE. pp. 783–794.

[31] Zhu, Q., Sun, Z., Liang, X., Xiong, Y., Zhang, L., 2020. Ocor: An overlapping-aware code retriever, in: 2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE. pp. 883–894.