Learning to Answer Programming Questions with Software Documentation through Social Context Embedding

Jing Lia,*, Aixin Suna, Zhenchang Xingb

^aSchool of Computer Science and Engineering, Nanyang Technological University, Singapore ^bCollege of Engineering and Computer Science, Australian National University, Australia

Abstract

Official software documentation provides a comprehensive overview of software usages, but not on specific programming tasks or use cases. Often there is a mismatch between the documentation and a question on a specific programming task because of different wordings. We observe from Stack Overflow that the best answers to programmers' questions often contain links to formal documentation. In this paper, we propose a novel deep-learning-to-answer framework, named QDLinker, for answering programming questions with software documentation. QDLinker learns from the large volume of discussions in community-based question answering site to bridge the semantic gap between programmers' questions and software documentation. Specifically, QDLinker learns question-documentation semantic representation from these question answering discussions with a four-layer neural network, and incorporates semantic and content features into a learning-to-rank schema. Our approach does not require manual feature engineering or external resources to infer the degree of relevance between a question and documentation. Through extensive experiments, results show that QDLinker effectively answers programming questions with direct links to software documentation. QDLinker significantly outperforms the baselines based on traditional retrieval models and Web search services dedicated for software documentation retrieval. The user study shows that QDLinker effectively bridges the semantic gap between the intent of a programming question and the content of software documentation.

^{*}Corresponding author

Email addresses: jli030@e.ntu.edu.sg (Jing Li), axsun@ntu.edu.sg (Aixin Sun), zhenchang.xing@anu.edu.au (Zhenchang Xing)

Keywords: Community-based question answering, Software documentation, Social context, Neural network

1. Introduction

For most programming languages and software packages, there exist comprehensive language specifications, Application Programming Interface (API) documentation, and tutorials. Such official documentation¹ provides information about functionality, structure, and parameters, but not on specific issues or specific usage scenarios [31, 42]. On the other hand, programmers often face very specific issues which are not explicitly stated in software documentation. For many such issues, software documentation does serve as a good reference for why the issues happen and how to address them. However, it is challenging to use a question as a keyword query to search for relevant software doc-

¹⁰ uments. This is because the software documentation and question are often in different wordings; one is for generic reference and the other is from a specific usage scenario in practice.

With the emergence of Web 2.0 in modern software development, the behavior of developers is changed, in relation to how they search for crowd-generated knowledge to fulfill their needs [21, 22, 25]. The mismatch between the needs of documentation consumers and the knowledge provided, leads to the overwhelming discussions accumulated at various Community-based Question Answering (CQA) websites such as Quora² and Stack Overflow³. In these discussions, the community users often refer to software documentation when answering programming questions. From Stack Over-

flow, we collected 45,288 best answers each contains at least one link to Java official documentation. Figure 1 plots the distribution of the number of links to Java documentation per best answer, which obeys a power-law distribution. It shows that 72.6% of best answers have exactly one link to Java documentation and fewer than 10% have more than three links. This distribution suggests that for many Java programming questions,

¹The term 'software documentation' refers to the collection of documents consisting of language specification, API documentation, and official tutorial.

²https://www.quora.com/

³http://stackoverflow.com/



Figure 1: Distribution of links to Java documentation per best answer, among 45, 288 best answers from Stack Overflow. The absolute numbers are plotted in log scale, and the percentages are plotted in bar chart.

- there exists a Java official document as a good reference to address the question. The large volume of discussions also create the '*semantic link*' between programmers' questions and software documentation, through the community of programmers, illustrated in Figure 2.
- Posting questions and waiting for answers from other programmers may take much time. The immediate question is: *can we answer a programmer's question by providing a link to the most relevant software documentation?* In this research, we aim to build an answering system where the questions are from programmers in natural language and the answers are the links to official documentation, illustrated in Figure 2. This system will provide convenience not only for documentation consumers but also the companies
- ³⁵ that provide technical support.

However, understanding programming questions to build an effective answering system is not trivial. First of all, mapping question-answer pairs into a discriminative feature space is a critical step. A widely adopted approach is to encode question-answer pairs using various features, *e.g.*, lexical, linguistic, and syntactic features [67, 37, 53,

⁴⁰ 61, 36]. These hand-crafted features may heavily depend on external resources at the loss of generality. Besides, many existing knowledge bases are about lexical knowledge or about open domain facts. A typical example is WordNet [29], a lexical knowledge



Figure 2: Overview of QDLinker. It directly links programmer's question to formal documentation through embedding the semantic context in CQA. Before QDLinker, the semantic links between questions and software documentation are established through the CQA community.

base for general English language, which may not be suitable to build answer systems for technical questions about programming. As shown by the analysis above, taking the ad-

- vantage of neural networks to learn semantic representation of question-documentation pair seems to be more appropriate for our task. Neural networks have been proved to be powerful tools in many fields, such as machine transliteration [7], computer vision [50], electromagnetic theory [18], wire coating analysis [30], and bioinformatics [40]. Note that, our task cannot be addressed by search engines for source code [13, 35]. Code
- so search system cannot well answer queries in natural language, especially when the queries do not contain any code snippets or API-like terms.

In this paper, we propose a novel *deep-learning-to-answer* framework named QDLinker, to answer programming questions with software documentation through social context embedding. *Social context* of a link to software documentation refers to the surrounding

- ⁵⁵ words of the link, when community users use it to answer questions in CQA. QDLinker embeds social contexts in a latent space, and uses a four-layer Deep Neural Network (DNN) to learn semantic representations of question-documentation pairs. The learned semantic representations and simple content features are then passed to a learning-torank schema to train a ranker. Compared to prior work on software text retrieval [67],
- our approach does not require manual feature engineering or hand-coded resources be-

yond the pre-trained word vectors. The architecture we proposed is beneficial not only to learn a ranker in training phase, but also to automatic feature extraction for the newcoming query-documentation pairs in online phase. Moreover, our approach takes into account documentation content and social context simultaneously, for its effectiveness

⁶⁵ in bridging the semantic gap between programming questions and software documentation.

We conducted extensive experiments on Stack Overflow dataset to evaluate the effectiveness of QDLinker. Empirical results show that QDLinker outperforms three baseline methods which are based on traditional retrieval models. Through a user study with 25 natural language queries collected from test dataset, we show that QDLinker significantly outperforms a commercial search engine. In short, our empirical results show that QDLinker can effectively bridge the semantic gap between questions and software documentation. In this paper, we make the following contributions:

• We propose QDLinker, a novel framework for answering programming questions

75

80

85

- with software documentation through social context embedding. It leverages the content in official sites and social contexts in CQA to learn semantic representations of question-documentation pairs and answers programming questions in natural language.
- We conduct a large-scale automatic evaluation, to evaluate the performance of QDLinker against three baseline methods. The empirical evaluation reveals that our approach can effectively answer Java technical questions against the traditional retrieval models.
- We conduct a user study to compare the software documentation retrieval performances of QDLinker and Google search. The results show that QDLinker significantly outperforms Google search in the retrieval task.

The remainder of this paper is organized as follows. Section 2 summarizes the related work. Section 3 details our approach QDLinker. Section 4 presents the empirical evaluation. Section 5 presents the user study. Finally, we conclude the paper in Section 6.

90 2. Related Work

95

Question Retrieval. Question retrieval has attracted much attention in recent years [4, 8, 16, 10]. Different retrieval models have been employed in the task, including the Okapi model [16], the translation model [63], the language model [8], and the vector space model [16, 17]. In addition, question category information has also been exploited for question retrieval [4]. Xue *et al.* [57] proposed a translation-based language model that combines the translation model and the language model for question retrieval. Yen *et al.* [61] developed a question classifier, which is trained to categorize the answer

type of a given question and instructs the context-ranking model to re-rank the passages retrieved from the initial retrievers.

- For the word mismatch problem among similar questions, existing solutions can be broadly grouped into three approaches to bridge the lexical gap. One approach is to use manual rules or templates. For example, Berger *et al.* [2] proposed a statistical lexicon correlation method to bridge the lexical chasm. The second approach is to use external knowledge databases such as Wikipedia and WordNet. The method by Zhou *et al.* [65]
- ¹⁰⁵ using semantic relations extracted from Wikipedia for question retrieval is an example. Burke *et al.* [3] proposed a model to rank frequently asked questions using combined similarities. The similarities are computed by conventional vector space models with semantic similarities based on WordNet. The third approach is to use deep representation. Zhou *et al.* [64, 66] proposed a neural network architecture to learn the semantic
- representations of question-answer pairs. Nassif *et al.* [33] presented a neural-based model with stacked bidirectional Long Short-Term Memory (LSTM) and Multi-Layer Perceptron (MLP) for similar question retrieval. Different from the prior studies, we aim to directly link questions to their relevant software documents, rather than retrieving similar questions.
- Answer Selection. Given a thread containing a question and a list of answers, many studies aim to automatically rank the answers according to their relevance to the question [49, 44, 51]. Sun *et al.* [49] used dependency relations between the matched question terms and the answer target as additional evidences to rank passages. Sakai *et al.* [44] proposed an approach to build answer selection system involving multiple an-

- swer assessors and graded-relevance information retrieval metrics. Yao *et al.* [60] proposed a family of algorithms to jointly detect the high-quality questions, and to help users to identify a useful answer that would gain much positive feedback from site users. Hou *et al.* [15] and Nicosia *et al.* [34] proposed automatic answer selection algorithms based on the position of the answer in the thread and the context of an answer in a thread,
- respectively. Instead of selecting answers, in our task, we attempt to distill the software documentation in answers.

CQA Semantic Representation. Most relevant to our work is the study on semantic representation of CQA. In recent years, deep neural networks have been used to learn higher-level semantic representations of question-answer pairs [46, 33, 47, 9, 23]. Tan *et*

- al. [52] developed hybrid models to match passage answers to questions accommodating their complex semantic relations. Severyn *et al.* [46] proposed a convolutional neural network architecture, which maps queries and documents to their distributed vectors, for reranking pairs of short texts. Yan *et al.* [58] proposed a deep neural network to learn how a query and its context are related to candidate reply. Nassif *et al.* [33] presented
- a neural-based model with stacked bidirectional LSTMs and MLP to learn semantic relatedness between questions and answers. Singh *et al.* [48] proposed a system using semantic keyword search in combination with traditional text search techniques to find similar questions with answers for unanswered questions.

Social Context. Yang *et al.* [59] investigated user preference and social contexts in ¹⁴⁰ point of interest (POI) recommendation. They developed a deep neural architecture that jointly learns the embeddings of users and POIs to predict both user preference and POIs. Bagci *et al.* [1] built a graph model of location-based social networks (LB-SNs) for personalized recommendations. This graph model took into account social contexts such as current social relations, personal preferences and current location of

the user. Li *et al.* [26] proposed WisLinker framework to recommend web resources using social contexts. Rohani *et al.* [43, 11] proposed an approach to address cold-start problem in academic social networks by incorporating social context features. Coined as an enhanced content-based algorithm using social networking (ECSN), the proposed algorithm considers the submitted ratings of faculty mates and friends besides user's



Figure 3: The architecture of our proposed approach. Offline phase aims to learn abstract representation for question-documentation pair, and to learn a ranker with the abstract representations. Given a new question, online phase aims to rank its candiate documents.

¹⁵⁰ own preferences. Rohani *et al.* [43] proposed an approach to solve cold-start problem by incorporating social context features.

In our proposed QDLinker, we learn the semantic representations of programming questions and the social contexts of software documentation. In our case, the social contexts, *i.e.*, how the APIs are used in different scenarios, cannot be obtained through the content of software documentation itself.

3. Deep Learning to Answer

In this section, we first give an overview of the proposed framework QDLinker, and then detail the core modules in QDLinker in Sections 3.2 - 3.4. The input to the framework, *i.e.*, the word embedding, is presented in Section 3.1.

160

155

As shown in Figure 3, the QDLinker framework consists of three core modules: candidate documentation generation, a four-layer neural network, and learning a ranker. Given a programming question in natural language, candidate documentation generation returns a small set of software documents which are considered relevant to the question. The DNN module learns the semantic representations of query-documentation

pairs in a latent space, and generates latent features for the ranker module. The features by DNN are fully automatic without human intervention. Nevertheless, the network does allow handcrafted features to be inserted in the join layer, illustrated in Figure 3. The learned features are then fed to a learning-to-rank schema to train a ranker, to pick up relatively relevant software documents among the candidates. Note that our approach

¹⁷⁰ cannot be formulated as an end-to-end problem (*i.e.*, directly minimizing a ranking cost function) because we need to automatically extract query-documentation representations for newcoming candidate documents in online phase. Thus, the DNN only serves as extracting the final representations of query-documentation pairs. Therefore, the architecture in Figure 3 is beneficial not only to learn a ranker in training phase, but also to

automatic feature extraction for newcoming query-documentation pairs in online phase.

3.1. Social Context Embedding

As shown in Figure 3, QDLinker is built on pre-trained word vectors, or word embeddings. Traditionally, language models represent each word as a feature vector using one-hot representation, where a vector element is 1 if the word is observed and

¹⁸⁰ 0 otherwise [12]. Recently, neural language models have been proposed to generate low-dimensional, distributed embeddings of words [6, 54]. These models take the advantage of word order in text documents and capture both syntactic and semantic relationships between words. Mikolov's continuous bag-of-words and skip-gram language models [27, 28] are among the most widely used models.

185

Stack Overflow is a destination with rich source of information about API usages and bug descriptions. Thus, in our implementation, we use the crowd-generated content on Stack Overflow to learn embeddings of words and links to software documentation.

As shown in Figure 4, a community user of Stack Overflow created a link to API documentation java.util.Collections.sort() in an answer. Figure 4 illustrates ¹⁹⁰ the training procedure with the skip-gram model when it reaches the link. This sentence creates a context for the link java.util.Collections.sort() through the surrounding words. We build two vocabularies: one for English words, and the other for links to software documentation. In simple words, each link to software documentation is treated as an ordinary term in the word sequence, and a word vector is learned for each link that is mentioned on Stack Overflow. Note that, we learn word embedding for each

link as a term, and the words in the anchor text of the link are not used in our training. We define that w_t is the only word on the input layer. N is the hidden layer size. V



Figure 4: Training word embedding example with the skip-gram model.

is the vocabulary size. *C* is the number of words in the context. $\mathbf{x} \in \mathbb{R}^{V}$ is the one-hot encoded vector for w_t , which means only one out of *V* units will be 1 and all other units are 0. The output of hidden layer can be written as

$$\boldsymbol{h} = \boldsymbol{W}^T \boldsymbol{x} = \boldsymbol{V}_{w_t}^T \tag{1}$$

where $\boldsymbol{W} \in \mathbb{R}^{V \times N}$ is the input-hidden weight matrix. \boldsymbol{V}_{w_t} is the vector representation of the input word w_t .

On the output layer, each output is computed using the hidden-output matrix:

$$p(w_{c,j} = w_{O,c} | w_t) = \frac{\exp(u_{c,j})}{\sum_{j'=1}^{V} \exp(u_{j'})}$$
(2)

where w_t is the input word. $w_{c,j}$ is the *j*-th word on the *c*-th panel of the output layer. $w_{O,c}$ is the actual *c*-th word in the output context word. $u_{c,j}$ is the net input of the *j*-th unit on the *c*-th panel of the output layer,

$$u_{c,j} = \boldsymbol{V'}_{w_j}^T \cdot \boldsymbol{h}, \text{ for } c = 1, 2, ..., C$$
(3)

where $\boldsymbol{V'}_{w_j}^T$ is the output vector of the *j*-th word in the vocabulary, w_j and $\boldsymbol{V'}_{w_j}^T$ is taken from a column of the hidden-output weight matrix, $\boldsymbol{W'}$.

When training the skip-gram model to predict C context words, the loss function is written as

$$E = -\log p(w_{O,1}, w_{O,2}, ..., w_{O,C} | w_t)$$

= $-\log \prod_{c=1}^{C} \frac{\exp(u_{c,j_c^*})}{\sum_{j'=1}^{V} \exp(u_{j'})}$ (4)



Figure 5: A 2D projection of embedding natural language terms and API documentation using PCA (API documentation in bold font and natural language terms in non-bold font).

where j_c^* is the index of the actual *c*-th output context word in the vocabulary.

Figure 5 illustrates a 2-D projection of vectors of natural language terms and API documentation using principal component analysis (PCA). In the embedding space, the vectors of terms and API documentation with the same intent have the shortest distance.

For example, the term "arraylist" is close to API documentation java.util.Arraylist. API documentation java.awt.Window.pack() is close to java.swing.JFrame. Observe that there are four clusters in Figure 5: java I/O, java.awt layout, java collection, and java compareto. For intra-cluster instances, they have similar programming function or purpose.

210 3.2. Candidate Documentation Generation

Given a programming question as a query, candidate generation selects a subset of software documents that are relevant to the question. We use three methods to select candidates.

Document Content. The content of software documentation reflects its relevance to a given query. In our implementation, we build a search engine for software documentation using Apache Lucene. Specifically, stopword removal and stemming are performed as preprocessing, and for each query, the search engine returns top 10 most relevant results based on the BM25 scoring function.

Local Context. Stack Overflow is a popular CQA site where developers ask questions and share knowledge about software development and maintenance. The discussions on Stack Overflow provide enriching context to mine usage scenarios of software documentation. When a software document appears in a discussion thread, its surround texts reflect its relevance to the question.

DEFINITION 1 (Local Context). If a software document is mentioned in a best answer, the texts of the question (title and body) and the best answer are regarded as the local context of the software document.

Local context is defined based on the consideration that the quality of best answer is better than other answers in the discussion thread, to avoid including too much noise. The body of the best answer is the immediate context where a software document is ²³⁰ mentioned. On the other hand, question title and body often describe the programming issues and reflect the relevance between the problem and the software documentation mentioned in its best answer.

Note that, each mention of a software document has its own local context. If a software document is mentioned multiple times, multiple local contexts are extracted.

We collect all local contexts of the mentioned software documents in our corpus. Given a query, we use Lucene to retrieve the most relevant local contexts, then pick the top 10 unique software documents as candidates.

Global Context. As aforementioned, a software document may be mentioned in multiple best answers and has multiple pieces of local contexts. For example, java.util.ArrayList was mentioned 906 times in our dataset.

DEFINITION 2 (Global Context). The global context of a software document is the collection of all its local contexts.

We build up a corpus through collecting global contexts for all software documents. Then we obtain the vector of a software document by social context embedding described in Section 3.1. Following [55], we use bag-of-words model to average out the vectors of the individual words in a query. Given a query, we retrieve top 10 software documents based on the cosine similarity between the average vector and software documentation vector.

3.3. Four-layer Deep Neural Network

Deep neural network with multiple layers has demonstrated its effectiveness in capturing semantical and higher-level discriminative information from the input data [24]. As shown in Figure 3, our DNN has four layers: *convolutional layer*, *join layer*, *hidden layer*, and *output layer*.

3.3.1. Convolutional Layer

255

250

For a natural language query with many words in a sentence, prior studies [62, 19] have shown that the simple bag-of-words model is unable to capture complex semantics of a sentence. Convolutional neural network can capture long-range dependencies and learn to correspond to the internal syntactic structure of sentences. Thus, we use one convolutional layer as the first layer in our approach.

Convolution Operation. Let X_d and X_q be the software documentation vector and query vector, respectively. Suppose that there are *s* words in the query and let $X_q^i \in \mathbb{R}^k$ be the *i*-th *k*-dimensional word vector corresponding to the *i*-th word in the query. More formally, the convolution operation * between two vectors $X_q \in \mathbb{R}^{sk}$ and $f_q \in \mathbb{R}^{mk}$ (called a filter of size *m*) results in a vector $c_q \in \mathbb{R}^{s-m+1}$ where each component is as follows:

$$c_q^j = (\boldsymbol{X}_q * \boldsymbol{f}_q)_j = \boldsymbol{f}_q^T \cdot \boldsymbol{X}_q^{[j:j+m-1]} + b_{qc}$$
(5)

where j = 1, ..., s - m + 1 and $\mathbf{X}_q^{[j:j+m-1]}$ represents the concatenation of word vectors $\mathbf{X}_q^j, \mathbf{X}_q^{j+1}, ..., \mathbf{X}_q^{j+m-1}, b_{qc} \in \mathbb{R}$ is a bias term. Thus, this filter \mathbf{f}_q is applied to each possible window of words in the query to produce a feature map:

$$\boldsymbol{c}_{q} = [c_{q}^{1}, c_{q}^{2}, \dots, c_{q}^{s-m+1}]$$
(6)

where $c_q \in \mathbb{R}^{s-m+1}$. Similarly, we can utilize filter f_d to produce documentation feature map c_d .

So far we have described the convolution layer with a single filter. Our model applies a set of filters that work in parallel to generate multiple feature maps. Let *n* be the number of filters. Given filters $F_q^{n \times mk}$ and $F_d^{n \times mk}$, the convolution operations produce two feature maps $C_q \in \mathbb{R}^{n \times (s-m+1)}$ and $C_d \in \mathbb{R}^{n \times (s-m+1)}$, respectively.

Activation Function. To allow the neural network to learn non-linear decision boundaries, each convolutional layer is followed by a non-linear activation function applied element-wise to the output of the convolution operations. Sigmoid, hyperbolic tangent *tanh*, and a rectifiled linear (*ReLU*) are among the most common choices for activation functions. In particular, it is reported that recified linear unit has significant benefits over sigmoid and *tanh* functions [32]. Thus, in our implementation, we use *ReLU* as

the activation function. The output of activation layer can be written as

$$\boldsymbol{A}_q = ReLU(\boldsymbol{C}_q) = max(0, \boldsymbol{C}_q) \tag{7}$$

$$\boldsymbol{A}_d = ReLU(\boldsymbol{C}_d) = max(0, \boldsymbol{C}_d)$$
(8)

where $\mathbf{A}_q \in \mathbb{R}^{n \times (s-m+1)}$ and $\mathbf{A}_d \in \mathbb{R}^{n \times (s-m+1)}$.

Pooling. The output from activation function is then passed to the pooling layer, whose goal is to aggregate the information and reduce the representation. As mentioned above, there are *n* filters. The pooling operation is applied on every filter. Taking the pooling of $A_q \in \mathbb{R}^{n \times (s-m+1)}$ as an example, the output of pooling $P_q \in \mathbb{R}^n$ can be written as

$$\boldsymbol{P}_{q} = \begin{bmatrix} pool\left(\boldsymbol{A}_{q}^{1}\right) \\ \dots \\ pool\left(\boldsymbol{A}_{q}^{n}\right) \end{bmatrix}$$
(9)

The pooling operation maps the feature map to a single value, formally: $pool(\mathbf{A}_q^i)$: $\mathbb{R}^{1 \times (s-m+1)} \rightarrow \mathbf{P}_q^i : \mathbb{R}$. There are a few common choices for the pool() operations: *average, max* and *L2-norm*. Average pooling was often used in the past but has recently fallen out of favor compared to the max pooling operation, which has been shown to work better in practice. In our approach, we use 1-*max* pooling strategy, which extracts a scalar with the maximum value for each feature map.

280 3.3.2. Join Layer

Inspired by [62, 46], we also add simple content features f_{cn} to our model. f_{cn} contains two word overlap features: word overlap count, and word overlap count weighted by IDF (inverse document frequency). Note that both features do not require any linguistic annotation or pre-processing. The output of join layer $X_{join} \in \mathbb{R}^{2n+2}$ can be expressed as follows:

$$\boldsymbol{X}_{join} = [\boldsymbol{P}_d; \boldsymbol{P}_q; \boldsymbol{f}_{cn}] \tag{10}$$

3.3.3. Hidden Layer

DNNs could use the intermediate layers to build up multiple layers of abstraction. These multiple layers of abstraction seem likely to give deep networks a compelling advantage in learning to solve complex pattern recognition problems [45]. In our architecture, the hidden layer is a fully-connected layer with parameters W_h and b. The output of hidden layer can be represented as

$$\boldsymbol{X}_{hidden} = ReLU(\boldsymbol{W}_h \cdot \boldsymbol{X}_{join} + \boldsymbol{b}_h)$$
(11)

3.3.4. Output Layer

The output of hidden layer X_{hidden} is passed to a fully connected softmax layer. It computes the probability distribution over the class labels:

$$p(y = j | \boldsymbol{X}_{hidden}; \boldsymbol{W}_s, \boldsymbol{b}_s) = softmax_j (\boldsymbol{W}_s \cdot \boldsymbol{X}_{hidden} + \boldsymbol{b}_s)$$
(12)

where \boldsymbol{W}_s and \boldsymbol{b}_s are the weight vector and the bias of softmax classifier, respectively. Our model is trained to minimize the cross-entropy cost function:

$$\mathscr{L} = -\log \prod_{i=1}^{N} p(y_i | \mathbf{X}_q^i, \mathbf{X}_d^i) + \lambda \| \boldsymbol{\theta} \|_2^2$$
(13)

where θ contains all parameters and we use *L2-norm* regularization.

$$\boldsymbol{\theta} = \left\{ \boldsymbol{F}_{q}, \boldsymbol{b}_{qc}, \boldsymbol{F}_{d}, \boldsymbol{b}_{dc}, \boldsymbol{W}_{h}, \boldsymbol{b}_{h}, \boldsymbol{W}_{s}, \boldsymbol{b}_{s} \right\}$$
(14)

In our problem setting, for a given question-documentation pair as an input instance, softmax layer outputs probabilities for two classification labels: positive and negative. Figure 6 shows an example of question-documentation pair extracted from discussions





on Stack Overflow. Together with the question, each link to documentation mentioned in the question's best answer forms a positive question-documentation pair instance.

For training the DNN, we use the links from the best answers of the training questions to form positive instances, and use randomly selected links to form negative instances. We use backpropogation algorithm to compute the gradients and use Adam update rule [20] to update the parameters of the network. To mitigate the overfitting issue, we augment the cost function with *L2-norm* regularization for the parameters of the network.

3.4. Learning a Ranker

295

In particular, X_{hidden} can be thought of as a final abstract representation of a querydocumentation pair, obtained by a series of transformations from the input layer through a series of layers. In our approach, we consider X_{hidden} as features to feed to a learningto-rank schema. The learning-to-rank schema can leverage multiple features for ranking and can automatically learn the optimal way of combining these features.

Our goal is to build a ranking model which facilitates each query q and its candidate list $D = \{d_1, d_2, ..., d_n\}$ to generate the optimal ranking. More formally, the task is to



Figure 7: The flow diagram of the proposed QDLinker. The middle part is to learn a model for extracting abstract representation for question-documentation pair. The left part aims to learn a ranker. The right part is to answer a new question with the learned models.

learn a scoring function F(q,d):

$$F(q,d) = \sum_{k=1}^{K} \omega_i \cdot \phi_i(q,d)$$
(15)

where each feature $\phi_i(q, d) \in \mathbf{X}_{hidden}$ measures a specific relationship between the query and a candidate software document. ω_i is the weight of the *i*-th feature (among the total *K* features), and is learned during the training process. In our task, the optimization procedure of learning-to-rank tries to find the scoring function that ranks the most relevant software document to the query at the top among all candidates. We train the ranking model using LambdaMART [56], a boosted tree version of LambdaRank [38], that won the Yahoo! Learning to Rank Challenge.

3.5. Summary and Complexity Analysis

To summarize, we present a flow diagram of QDLinker in Figure 7 and the pseudocode in Algorithm 1 and Algorithm 2. Observe from Figure 7, our framework contains both offline phase and online phase. The offline phase first learns abstract representation of question-documentation pair and then learns a ranker based on the positive and negative instances. Here an instance is a question-documentation pair. The offline

Algorithm 1: Pseudocode for offline phase of the proposed QDLinker

| I | Input: N training instances of question-document pairs | | | | | | | | |
|------|---|--|--|--|--|--|--|--|--|
| 0 | Output: Feature vector function $V_{\theta}(q,d)$ and ranking function <i>F</i> | | | | | | | | |
| / | // Phase 1: Learning abstract representation | | | | | | | | |
| 1 II | 1 Initialize all parameters θ ; | | | | | | | | |
| 2 fc | e foreach <i>epoch in epoch_{max}</i> do // iterate through epoches | | | | | | | | |
| 3 | Sample a minin-batch from N pairs; | | | | | | | | |
| 4 | Clear gradients $d\theta \leftarrow 0$; | | | | | | | | |
| 5 | Computing \mathscr{L} based on Equation (13); | | | | | | | | |
| 6 | Update $	heta \leftarrow 	heta - rac{\partial \mathscr{L}}{\partial 	heta} \cdot lr$; | | | | | | | | |
| 7 | return $V_{	heta}(q,d)$; // return feature vector function for q-d pair | | | | | | | | |
| / | / Phase 2: Learning a ranker | | | | | | | | |
| 8 S | et number of trees M , number of leaves per tree L ; | | | | | | | | |
| 9 fc | oreach m in M do // iterate trees | | | | | | | | |
| 10 | foreach n in N do // iterate through training pairs | | | | | | | | |
| 11 | Calculating feature vector for current <i>n</i> via $V_{\theta}(q, d)$ in Phase 1; | | | | | | | | |
| 12 | Calculating the λ -gradients for each q-d pair; // more detials[56] | | | | | | | | |
| 13 | Calculating the second-order derivative using the λ -gradients; | | | | | | | | |
| 14 | Building a regreesion tree with <i>L</i> terminal nodes; | | | | | | | | |
| 15 | Update F function; | | | | | | | | |
| 16 | return F; // Final ranking function F | | | | | | | | |

phase is listed in Algorithm 1. Given a new question, the online phase first generates the k candidate software documents based on Section 3.2. Then these documents are ranked by the learned models in offline phase. Accordingly, the online phase is listed in Algorithm 2.

315

320

We first detail the parameter size of our framework, then deduce the time and space complexity. Equations (14) and (15) show all parameters that QDLinker would learn. Now we consider one query-document pair as shown in Figure 3. For convolutional layer, QDLinker has $2n \times mk + 2n$ parameters, where *n* is the number of filters and *m*

- is the window size of each filter, k is the dimension of word embedding, and number 2 indicates query channel and document channel. For join layer, there is no parameter. For hidden layer, there are $(2n+2) \times h + h$ parameters, where h is the number of neruons in hidden layer. For output layer, there are $h \times 2 + 2$ parameters. For ranker layer, there
- are *h* weights because QDLinker uses the X_{hidden} as the final abstract representation. Given a new question, we assume that QDLinker generates κ candidate software

Algorithm 2: Pseudocode for online phase of the proposed QDLinker

| | Input: A new question q_{new} |
|---|---|
| | Output: top- <i>k</i> software documents |
| 1 | Retrieving candidate documents for q_{new} based on Section 3.2; |
| 2 | Extracting features for these question-document pairs using $V_{\theta}(q, d)$ in Phase 1 of Algorithm 1; |
| 3 | Ranking these candidate documents using F in Phase 2 of Algorithm 1; |

4 return top-k software documents; // Answers to the issued question

documents. Now considering a question-documentation pair, the total time and space complexity of convolutional layer are both $O(\alpha \cdot m^2 \cdot \beta)$ [5, 14], where α and β are the number of input nodes and the number of output nodes respectively, *m* is the window size of each filter. Hidden layer and output layer are fully connected layers. The time and space complexity of this linear projection are both $O(\alpha\beta)$. For the ranker, the time complexity is $O(\kappa^2)$ [56]. Thus, the time complexity of QDLinker is $O(\kappa\alpha m^2\beta + \kappa^2)$,

and the space complexity is $O(\alpha m^2 \beta)$.

4. Empirical Evaluation

330

We now evaluate the effectiveness of QDLinker by measuring its accuracy on linking questions on Stack Overflow to software documentation. Our evaluation assumes that the software documentation mentioned in a question's best answer is the most relevant to the question.

4.1. Experimental Setting

- Data Collection. In our evaluation, we focus on Java software documentation which consists of Java Standard Edition API documentation, Java tutorials, and language specifications. Usually, a programming question is expressed in natural language thus it is similar to the discussions on Stack Overflow. We therefore use the data collected from Stack Overflow in our experiments.
- We extract discussion threads from the datadump archive⁴ that satisfy the following criteria: (i) The score of question is greater than 0. This condition guarantees that at

⁴https://archive.org/details/stackexchange

| Data | #Discussion threads |
|-----------------|---------------------|
| Word embeddings | 24,217 |
| Train | 10,649 |
| Development | 1,000 |
| Test | 1,693 |

Table 1: Summary of the number of threads used in each task.

least one developer has voted the question to be a 'useful question'. (ii) The question has an answer which is accepted as the best answer, and the score of the best answer is greater than 0, and (iii) The best answer must contain at least one link to the above listed

Java documentation. Based on the above criteria, we collect 30,272 discussion threads from the data dump released on August 2015. We randomly select 24,217 discussion threads (account for 80%) as training data and the remaining 6,055 threads (account for 20%) as test data.

Model Training. We learn word embeddings from the training data, *i.e.*, the 24,217
³⁵⁵ discussion threads. For each discussion thread, we extract text from question title, question body, and all answers whose scores are greater than 0. We use the skip-gram model implemented in word2vec⁵. The context window size is set to 10 and the minimal word frequency is 5. Recall that each link to a software document is also treated as a word (or term, see Figure 4). Based on this condition, we have 1,520 distinct links to Java
³⁶⁰ documentation in the training data.

Next is to train the DNN and the ranker. Note that some discussion threads cannot be used to extract query-documentation pairs for training DNN because the links to software documentation in best answers are filtered out when the minimal word frequency is set to 5. Finally, we have 10,649 discussion threads for training DNN and the ranker,

1,000 discussion threads used for development set, and 1,693 discussion threads used for test. Table 1 summarizes the dataset in our experiments.

We empirically set the hyperparameters based on the development set. The number

⁵https://code.google.com/archive/p/word2vec/

of filters in convolutional layer is 64, and the size of filter is set to 2. The size of hidden layer is set to 64. The dimensionality of pre-trained word vectors is 200. *L2-norm* term is set to 1e-5 and the learning rate is 1e-3.

Performance Measures. We use the following five performance metrics in our evaluation:

- Precision at k, $P@k = \frac{|D_k \cap D_g|}{k}$, is the fraction of relevant documentation links to the query question among the top k ranked results. D_k denotes the set of top-k ranked links to software documentation and D_g is the set of ground-truth links (*i.e.*, links to software documentation in the question's best answer).
- Recall at *k*, $R@k = \frac{|D_k \cap D_g|}{|D_g|}$ is the fraction of ground-truth links in the top-*k* results.
- Hit rate at k, denoted by HR@k, is 1 if $|D_k \cap D_g| > 0$ and 0 otherwise.
- Mean average precision, *MAP*, is the mean of average precision (AP) over a set of test queries.
 - Mean reciprocal rank, $MRR(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{rank_j}$, is the average reciprocal rank of the results over a set of test queries Q. In the equation, $rank_j$ denotes the rank position of the first relevant document for the *j*-th test query.
- **Baseline Methods.** In order to validate the effectiveness of the proposed method, we evaluate the following three baseline methods in our experiments.
 - OfficialCn: This is the baseline model which selects candidate software documentation by content (see Section 3.2). BM25 [41] scoring function is used to rank the candidates.
- LocalCx: This baseline ranks candidates based on local context (see Section 3.2). This model indexes the local contexts in discussion threads and retrieves candidates using BM25 scoring function.

375

380

Table 2: Performance (P@k, R@k, MAP and MRR) for different methods. The best performance is highlighted in bold face. \dagger indicates that the differences between the result of QDLinker and other models are statistically significant with p < 0.05 under *t*-test.

| Method | P@1 | P@2 | P@5 | P@10 | R@1 | R@2 | R@5 | R@10 | MAP | MRR |
|------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| OfficialCn | 0.1347 | 0.1087 | 0.0749 | 0.0511 | 0.1147 | 0.1797 | 0.3053 | 0.4108 | 0.2234 | 0.2267 |
| LocalCx | 0.1630 | 0.1337 | 0.0932 | 0.0657 | 0.1401 | 0.2253 | 0.3876 | 0.5361 | 0.2896 | 0.2956 |
| GlobalCx | 0.1536 | 0.1234 | 0.0875 | 0.0628 | 0.1300 | 0.2002 | 0.3548 | 0.5005 | 0.2564 | 0.2614 |
| QDLinker | 0.1875† | 0.1576† | 0.1272† | 0.0919† | 0.1708† | 0.2734† | 0.5012† | 0.6847† | 0.3461† | 0.3584† |

Table 3: Performance(HR@k) for different methods.

| Method | HR@1 | HR@2 | HR@5 | HR@10 |
|------------|---------|---------|---------|---------|
| OfficialCn | 0.1347 | 0.2069 | 0.3491 | 0.4663 |
| LocalCx | 0.1630 | 0.2631 | 0.4341 | 0.5914 |
| GlobalCx | 0.1536 | 0.2383 | 0.4109 | 0.5640 |
| QDLinker | 0.1875† | 0.3057† | 0.5828† | 0.8128† |

GlobalCx: This baseline ranks candidates based on global context (see Section 3.2). This model represents natural language words and software documentation as vectors in shared embedding space [27].

395

4.2. Performance Comparison

Table 2 and Table 3 report the linking performance by different methods on all evaluation metrics. From the results, QDLinker significantly outperforms all baselines on all metrics. The improvement is statistically significant based on *t*-test with p < 0.05.

- ⁴⁰⁰ Observe that P@10 is very small for all methods including QDLinker. On Stack Overflow, more than 70% of the best answers contain only one link to software documentation (reported in Section 1, Figure 1). As the result, $|D_{10} \cap D_g| = 1$ in most cases. Thus, the ideal value of P@10 is slightly above 0.1, QDLinker achieves a very good result of 0.0919 in this sense. On R@k measure, QDLinker significantly outperforms the other
- baselines for all k values (k = 1, 2, 5, 10). It is worth noting that QDLinker achieves the highest recall (0.6847) when k = 10. In terms of *MAP* measure, QDLinker outperforms the three baseline methods OfficialCn, LocalCx and GlobalCx by 34.98%, 19.51% and 54.93%, respectively. On *MRR* measure, QDLinker outperforms the three baselines by 37.11%, 21.24% and 58.09%, respectively. Similar observations hold on hit
- ⁴¹⁰ rate measure HR@k, reported in Table 3.

Table 4: Results on additional content features.

| Content features | MAP | MRR | | |
|---------------------------|-----------------|----------------------------|--|--|
| Without f_{cn} features | 0.3054 | 0.3102 | | |
| With f_{cn} features | 0.3461(†13.32%) | $0.3584(\uparrow 15.53\%)$ | | |

Observe from the results that the content-based method OfficialCn delivers the worst performance on all measures. This implies that content-based approach cannot bridge developers' intent and content of software documentation. This observation is consistent with our description in Section 1 that software documentation is prepared to give comprehensive coverage without targeting on specific problems, while the programming questions are encountered in specific programming tasks. Therefore, it is essential to utilize the social context available on Stack Overflow to bridge the semantic gap between programmers' questions and software documentation.

4.3. Impacts of Factors on Performance

420 4.3.1. Impact of content features

In our approach, f_{cn} consists of word overlap count and word overlap count weighted by IDF value. Table 4 shows the performance and improvement of considering additional content features f_{cn} . More specifically, considering f_{cn} improves *MAP* by 13.32%, and *MRR* by 15.53%.

- There are two aspects resulting in the improvement. First, as described above, input of our approach are the pre-trained word vectors. In fact, the word dictionary of our dataset may not cover all the words in English language. The overlap features can provide supplementary information in our approach. Additionally, one of the weaknesses of approaches relying on distributed word vectors is their inability to deal with numbers
- and proper nouns [62, 46]. But when developers issue natural language queries, most of the questions are of type "what", "when", "who" that are looking for answers containing numbers or proper nouns. Thus, the model with f_{cn} features outperforms the one without f_{cn} features on *MAP* and *MRR*.



Figure 8: Impact of dimensionality of word embedding.

4.3.2. Dimensionality of word embedding

435

440

445

QDLinker takes in pre-trained word vectors in the input layer and feeds into the convolutional layer. We vary the dimensionality of word embedding and evaluate its impact on *MAP* and *MRR*. Figure 8 reports the performance of QDLinker using word embedding in different dimensions (50, 100, 200, 300, 400, 500 and 600). The results indicate that the dimensionality of word embedding has very marginal impact on the performance.

One possible reason is that the distributed word vectors varying different dimensions contain enough latent information for building a ranker in our dataset. Thus, when training a ranker, our approach is stable for different dimensions of word embedding.

4.3.3. Impact of layer sizes

As shown in Figure 3, out architecture needs to set the number of neurons (*i.e.*, layer size) in convolutional layer and hidden layer. Figure 9 shows the impact of setting different convolutional layer sizes and hidden layer sizes, on the test set. The comparison is based on 200 dimensions of word vectors.

We observe that the performance of ranker greatly depends on the combination of convolutional layer size and hidden layer size. In our dataset, we obtain the best performance when the convolutional layer size and hidden layer size are both 64.

When the combination layer sizes are small, the *MAP* is around 0.2 only, much lower than the best performance. Too few neurons in the convolutional and hidden layers will result in underfitting, as the neurons cannot capture enough signals to model complex



Figure 9: Performance in MAP with different layer sizes.

data. However, larger combination layer sizes does not lead to better ranker performance either. In addition, a large number of neurons in the convolutional and hidden layers increase the training time.

In summary, dimensionality of word vectors has very marginal impact on the performance. However, the sizes of convolutional and hidden layers have significant impact on the performance of our model.

5. User Study

To the best of our knowledge, there is no existing work on answering programming questions in natural language. Commercial search engines, *e.g.*, Google and Bing, are tools for daily use in software development. It naturally motivates us to compare the returned results with such search engines. If we can improve the performance of search results on the search engines, it will provide convenience not only for developers but also for the companies that provide documentation support.

In the previous set of experiments, we consider the software documentation mentioned in the best answers as the ground truth to questions. This assumption may ignore ⁴⁷⁰ the other retrieved software documentation which is relevant to the question but is not mentioned in the best answers. That is, although limiting documentation mention in best answer is a good criterion to control the quality of ground truth, the criterion may exclude the relevant results from our evaluation. In this section, we perform a user study to manually evaluate performance of QDLinker against Web search services.

475 5.1. Evaluation Setup

From the test dataset, we randomly select 25 discussion threads and query QDLinker using their questions. The 25 questions are listed in Table 5. For each question, we also use Google search engine to retrieve a list of software documentation. Because we focus on official documentation in this study, we restrict the retrieved results by

- Google by setting the "site" parameter in the search. For example, the second query in Table 5 is extended as"*XML string parsing in Java site:docs.oracle.com/javase*" using Google search engine in August, 2016. Note that, all the three types of Java documentation (language specification, API documentation, and tutorial) are under the same site: docs.oracle.com/javase.
- To measure the performance of QDLinker and Google, we use three metrics [13, 39]. *FR* is the rank of the first relevant result, as most users scan the results from top to bottom. The smaller the number of *FR*, the better the performance. The P@5 denotes the precision of the top 5 ranked results. Note that, the judgements of relevance are manually labeled by our annotators. Similarly, the P@10 is the precision of the top 10 ranked results.

We recruited two developers to manually annotate the two set of results from Google and QDLinker, respectively. Each link to software documentation in result list was marked relevant or irrelevant, indicating whether the developer considered this software documentation is relevant to the question. The annotation was done individually by the two developers and for inconsistent judgments, the two developers reached a consensus through discussion.

5.2. Evaluation Results

495

Table 5 shows the performance comparison of Google search and QDLinker. In particular, the symbol "-" in the second column indicates that there is no relevant soft-

Table 5: Human evaluation for Google search and QDLinker (*BA*: the number of ground truth links in best answer. *FR*: the rank of the first relevant documentation. P@5 and P@10: precision of the first 5 and 10 results. "A_" indicates Java API documentatins. "S_" indicates Java language specification. "T_" indicates Java tutorial. The boldface documentation indicate the ground truth documentation in best answers. † indicates the differences between QDLinker and Google are significant with p < 0.05 under *t*-test).

| | | Google Search | | QDLinker | | er | | |
|--|-----|---------------|--------|----------|-------|--------|-------|---|
| Query | BA | FR | P@5 | P@10 | FR | P@5 | P@10 | Top 3 relevant documents by QDLinker |
| 1: java switch error when simplifying code | 1 | 7 | 0 | 0.2 | 1 | 0.8 | 0.5 | S_Chapter 16. Definite Assignment; T_Branching Statements; S_Chapter 14.11. The switch Statement |
| 2: XML string parsing in Java | 1 | 1 | 1 | 0.9 | 4 | 0.4 | 0.3 | $\label{eq:converter} {\bf A_javax.xml.parsers.DocumentBuilder.parse(); {\bf A_javax.xml.bind.DatatypeConverter; {\bf A_javax.xml.transform.Transformer} = 0.0000000000000000000000000000000000$ |
| 3: PLAF can't change button color | 2 | 1 | 0.4 | 0.5 | 1 | 0.4 | 0.3 | A_javax.swing.UIManager; T_How to Set the Look and Feel; T_How to Use Color Choosers |
| 4: match generics with mockito | 1 | 1 | 0.4 | 0.5 | 1 | 0.6 | 0.7 | T_Type Erasure; T_Erasure of Generic Types; T_Erasure of Generic Methods |
| 5: reinitialise transient variable | 1 | | 0 | 0 | 1 | 0.6 | 0.7 | A_java.io.Serializable; S_Chapter 8.3.2. Initialization of Fields; S_Chapter 14.14. The for Statement |
| 6: write and read multiple byte[] in file | 1 | 2 | 0.2 | 0.2 | 1 | 0.6 | 0.5 | $\label{eq:alpha} A_java.io.FileInputStream.read(); A_java.io.FileOutputStream.write(); A_java.io.RandomAccessFileStream.read(); A_java.io.FileOutputStream.write(); A_java.io.RandomAccessFileStream.write(); A_java.io.FileOutputStream.write(); A_$ |
| 7: equality of boxed boolean | 1 | 1 | 0.4 | 0.5 | 1 | 0.6 | 0.5 | S_Chapter 5.1.7. Boxing Conversion; A_java.util.Arrays.equals(); A_java.lang.Object.equals() |
| 8: resetting and copying two dimensional arrays | 1 | - | 0 | 0 | 1 | 0.4 | 0.4 | A_java.lang.System.arraycopy(); T_Arrays; A_java.util.Arrays.copyOf() |
| 9: java standard on result of casting a double to an int | 2 | 1 | 0.8 | 0.6 | 1 | 0.8 | 0.7 | S_Chapter 5.1.3. Narrowing Primitive Conversion; A_java.lang.Number; A_java.lang.Double |
| 10: registering and using a custom java.net.URL protocol | 5 | 3 | 0.2 | 0.1 | 1 | 0.6 | 0.8 | A_java.net.URLConnection; A_java.net.URI; A_java.net.HttpURLConnection |
| 11: why is the protected method not visible | 1 | 1 | 0.2 | 0.1 | 3 | 0.4 | 0.5 | S_Chapter 5.1.3. Narrowing Primitive Conversion; S_Chapter 5.1.2. Widening Primitive Conversion T_Controlling Access to |
| | | | | | | | | Members of a Class; |
| 12: java date formatting ParseException | 1 | 1 | 0.8 | 0.6 | 1 | 0.6 | 0.7 | $\label{eq:alpha} A_java.text.SimpleDateFormat; A_java.text.DateFormat.parse(); A_java.text.DateFormat.format() \\$ |
| 13: java executor with no ability to queue tasks | 2 | - | 0 | 0 | 1 | 0.8 | 0.8 | A_java.util.concurrent.Executors; A_java.util.concurrent.ThreadPoolExecutor; A_java.util.concurrent.ExecutorService |
| 14: how to replace a jPanel based on user clicks | 1 | 2 | 0.6 | 0.5 | 2 | 0.4 | 0.6 | T_How to Use CardLayout; T_How to Use the Focus Subsystem; T_How to Write a Mouse Listener |
| 15: ambiguous varargs method call compilation error | 1 | - | 0 | 0 | 2 | 0.4 | 0.5 | S_Chapter 15.12.2. Compile-Time Step 2: Determine Method Signature; A_java.lang.invoke.MethodHandle; S_Chapter |
| | | | | | | | | 6.6.1. Determining Accessibility |
| 16: why is there no generic type information at run rime | 2 | 1 | 0.8 | 0.7 | 1 | 0.6 | 0.6 | T_Generic Methods; T_Why Use Generics?; T_Erasure of Generic Types |
| 17: raster format exception (Y+height) | 1 | 4 | 0.2 | 0.3 | 1 | 0.4 | 0.4 | A_java.awt.image.BufferedImage; A_javax.imageio.ImageIO; T_Lesson: Working with Images |
| 18: rendering combo boxes in a JTable | 2 | 1 | 0.6 | 0.6 | 1 | 0.8 | 0.7 | T_How to Use Tables; T_How to Write an Item Listener; T_How to Use Combo Boxes |
| 19: java modifying a class directly, null reference | 1 | - | 0 | 0 | 1 | 0.6 | 0.5 | T_Anonymous Classes; S_Chapter 12.4. Initialization of Classes and Interfaces; A_java.lang.NullPointerException |
| 20: windows azure date format to java date | 1 | - | 0 | 0 | 1 | 0.8 | 0.7 | A_java.text.SimpleDateFormat; A_java.time.format.DateTimeFormatter; A_java.text.DateFormat |
| 21: reading arraylist from a .txt file | 1 | 5 | 0 | 0.1 | 2 | 0.6 | 0.6 | A_java.nio.file.Files.readAllLines(); A_java.io.FileInputStream; A_java.util.Scanner.nextLine(); |
| 22: close connection and statement finally | 1 | 1 | 0.8 | 0.4 | 1 | 0.8 | 0.6 | T_The try-with-resources Statement; A_java.sql.Statement; T_Using Transactions |
| 23: when are java temporary files deleted | 2 | 1 | 0.4 | 0.2 | 1 | 0.4 | 0.5 | $\label{eq:linear} A_java.io.File.createTempFile(); A_java.io.File.deleteOnExit(); A_java.nio.file.Files.createTempDirectory() \\ \label{eq:linear}$ |
| 24: shutdown application gracefully upon power loss | 1 | - | 0 | 0 | 1 | 0.6 | 0.6 | $\label{eq:lang_relation} A_java.lang.Runtime.addShutdownHook(); A_java.util.concurrent.ExecutorService.shutdownNow() A_java.util.Timer = 100000000000000000000000000000000000$ |
| 25: get single bytes from multi-byte variable | 1 | 1 | 0.2 | 0.1 | 1 | 0.4 | 0.5 | T_Primitive Data Types;A_java.nio.ByteBuffer ; T_Variables |
| average | 1.4 | > 1.94 | 4 0.32 | 0.284 | 1.32† | 0.576† | 0.568 | |

ware documentation returned by Google search in the query. The last row shows the average performance on the three metrics.

Compared with Google search, QDLinker achieves better performance on FR, P@5 and P@10. In most cases (20 out of the 25 queries), QDLinker is able to recommend relevant software documentation at the first position in the result list. The differences between these two approaches in terms of the three metrics are statistically significant at p < 0.05. That is, QDLinker provides more relevant software documentation in top 10 results than Google search in our user study.

The last column shows web page titles of the top 3 ranked relevant documents by QDLinker, which consist of Java API documentation (marked by $A_{)}$, Java language specifications (marked by $S_{)}$ and Java tutorials (marked by $T_{)}$. For example, in

Queries 1 and 2, "S_Chapter 16. Definite Assignment"⁶, "A_javax.xml.parsers. Docu-

⁶https://docs.oracle.com/javase/specs/jls/se8/html/jls-16.html

mentBuilder. parse()"⁷, "*T_Branching Statements*"⁸ represent a document in language specification, API documentation, and tutorial, respectively. Compared with Google search, we make following observations:

- QDLinker can bridge the semantic gap between question and software documentation. For example, query 5 is "reinitialise transient variable", and there is no Java documentation which contains all the three keywords. Google search cannot return relevant results in this query. Likewise, the state-of-the-art API usage miner [13] cannot return any API sequences based on code corpus from Github.
 We manually check our training dataset and find that some community users have implemented the task using the class "java.io.Serializable" and method "readObject" on Stack Overflow. Thus, QDLinker can effectively answer this question because it takes into account the content and context of software documentation simultaneously.
 - QDLinker can effectively answer complex and bug-like queries. For instance, query 17 "raster format exception (Y+height)" and query 19 "java modifying a class directly, null reference" are related to program exceptions. Poorer results were obtained from Google for such kind of queries. On the contrary, QDLinker provided high quality results for such kind of queries, and an example is the API documentation java.awt.image.BufferedImage for query 17.

525

530

535

• QDLinker can effectively answer programming questions which are in specific usage scenarios. For instance, query 18 "*rendering combo boxes in a JTable*" is about usage of *combo boxes* in the scenario of "*JTable*" and query 13 "*java executor with no ability to queue tasks*" is about of *Java executor* in the scenario of "*queue tasks*". The official software documentation does not serve as good reference for these queries in specific usage scenarios, while QDLinker can provide high quality software documentation for these queries. For example,

⁷https://docs.oracle.com/javase/8/docs/api/javax/xml/parsers/DocumentBuilder.html\
#parse-java.io.File-

⁸https://docs.oracle.com/javase/tutorial/java/nutsandbolts/branch.html

java.util.concurrent.ThreadPoolExecutor is a high-quality API document for query 13.

540 6. Conclusion

Developers often encounter questions in specific programming tasks. Although programming languages and software packages are well supported by formal documentation, the documentation aims at comprehensive coverage and not on specific tasks. The semantic gap between the developers' questions and software documentation makes it

- ⁵⁴⁵ difficult for developers to search for the most relevant documentation. Utilizing the social context available on Stack Overflow, we built QDLinker to bridge the gap between the questions and documentation. Given a programming question, QDLinker returns the links to the most relevant documentation. The semantic features between questions and software documentation in QDLinker are learned through a four-layer deep neural
- network. Together with content features, the learned features are fed to a learning-torank schema for ranking the most relevant software documentation at top position. Using real questions from Stack Overflow, we show that QDLinker effectively locates the most relevant software documentation to questions, and its performance significantly outperforms baseline methods.
- The proposed QDLinker framework may benefit other software engineering problems. First, considering the ability of bridging the semantic gap between programming questions and software documentation, QDLinker could improve official software documentation with the information in questions and answers. Second, current code search does not support natural language. QDLinker can be integrated in code search en-
- gines to improve code search performance. Third, this work opens several interesting directions for future work with regard to automatic conversation between humans and computers. In the future, we will explore the applications of QDLinker to these problems.

References

- [1] Bagci, H., & Karagoz, P. (2016). Context-aware location recommendation by using a random walk-based approach. *Knowledge and Information Systems*, 47, 241–260.
 - [2] Berger, A., Caruana, R., Cohn, D., Freitag, D., & Mittal, V. (2000). Bridging the lexical chasm: statistical approaches to answer-finding. In *SIGIR* (pp. 192–199).
- ⁵⁷⁰ [3] Burke, R. D., Hammond, K. J., Kulyukin, V., Lytinen, S. L., Tomuro, N., & Schoenberg, S. (1997). Question answering from frequently asked question files: Experiences with the faq finder system. *AI magazine*, *18*, 57.
 - [4] Cao, X., Cong, G., Cui, B., & Jensen, C. S. (2010). A generalized framework of exploring category information for question retrieval in community question answer archives. In WWW (pp. 201–210).
- 575

- [5] Cheng, Y., Yu, F. X., Feris, R. S., Kumar, S., Choudhary, A., & Chang, S.-F. (2015). An exploration of parameter redundancy in deep networks with circulant projections. In *ICCV* (pp. 2857–2865).
- [6] Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, *12*, 2493–2537.
 - [7] Deselaers, T., Hasan, S., Bender, O., & Ney, H. (2009). A deep learning approach to machine transliteration. In *Proceedings of the Fourth Workshop on Statistical Machine Translation* (pp. 233–241).
- [8] Duan, H., Cao, Y., Lin, C.-Y., & Yu, Y. (2008). Searching questions by identifying question topic and question focus. In ACL (pp. 156–164).
 - [9] Er, M. J., Zhang, Y., Wang, N., & Pratama, M. (2016). Attention pooling-based convolutional neural network for sentence modelling. *Information Sciences*, 373, 388–403.

- ⁵⁹⁰ [10] Figueroa, A., & Neumann, G. (2016). Context-aware semantic classification of search queries for browsing community question–answering archives. *Knowledge-Based Systems*, 96, 1–13.
 - [11] Gani, A., Siddiqa, A., Shamshirband, S., & Hanum, F. (2016). A survey on indexing techniques for big data: taxonomy and performance evaluation. *Knowledge* and Information Systems, 46, 241–284.
- 595

- [12] Grbovic, M., Djuric, N., Radosavljevic, V., Silvestri, F., & Bhamidipati, N. (2015).
 Context-and content-aware embeddings for query rewriting in sponsored search. In *SIGIR* (pp. 383–392).
- [13] Gu, X., Zhang, H., Zhang, D., & Kim, S. (2016). Deep api learning. In *FSE* (pp. 631–642). ACM.
- [14] He, K., & Sun, J. (2015). Convolutional neural networks at constrained time cost. In *CVPR* (pp. 5353–5360).
- [15] Hou, Y., Tan, C., Wang, X., Zhang, Y., Xu, J., & Chen, Q. (2015). Hitszicrc: Exploiting classification approach for answer selection in community question answering. In *SemEval* (pp. 196–202). volume 15.
- [16] Jeon, J., Croft, W. B., & Lee, J. H. (2005). Finding similar questions in large question and answer archives. In *CIKM* (pp. 84–90).
- [17] Ji, Z., Xu, F., Wang, B., & He, B. (2012). Question-answer topic model for question retrieval in community question answering. In *CIKM* (pp. 2471–2474).
- ⁶¹⁰ [18] Khan, J. A., Raja, M. A. Z., Rashidi, M. M., Syam, M. I., & Wazwaz, A. M. (2015). Nature-inspired computing approach for solving non-linear singular emden–fowler problem arising in electromagnetic theory. *Connection Science*, 27, 377–396.
- [19] Kim, Y. (2014). Convolutional neural networks for sentence classification. arXiv
 preprint arXiv:1408.5882, .

- [20] Kingma, D., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, .
- [21] Ko, A. J., DeLine, R., & Venolia, G. (2007). Information needs in collocated software development teams. In *ICSE* (pp. 344–353).
- ⁶²⁰ [22] Ko, A. J., Myers, B. A., Coblenz, M. J., & Aung, H. H. (2006). An exploratory study of how developers seek, relate, and collect relevant information during software maintenance tasks. *IEEE Transactions on software engineering*, *32*, 971–987.
 - [23] Kokkinos, Y., & Margaritis, K. G. (2015). Topology and simulations of a hierarchical markovian radial basis function neural network classifier. *Information Sciences*, 294, 612–627.
 - [24] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521, 436– 444.
 - [25] Li, J., Bao, L., Xing, Z., Wang, X., & Zhou, B. (2016). Bpminer: mining develop-

- ers' behavior patterns from screen-captured task videos. In SAC (pp. 1371–1377). ACM.
- [26] Li, J., Xing, Z., Ye, D., & Zhao, X. (2016). From discussion to wisdom: web resource recommendation for hyperlinks in stack overflow. In SAC (pp. 1127– 1133).
- 635 [27] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, .
 - [28] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *NIPS* (pp. 3111–3119).
- [29] Miller, G. A. (1995). Wordnet: a lexical database for english. *Communications of the ACM*, 38, 39–41.

- [30] Munir, A., Manzar, M. A., Khan, N. A., & Raja, M. A. Z. (). Intelligent computing approach to analyze the dynamics of wire coating with oldroyd 8-constant fluid. *Neural Computing and Applications*, (pp. 1–25).
- ⁶⁴⁵ [31] Nadi, S., Krüger, S., Mezini, M., & Bodden, E. (2016). Jumping through hoops: why do java developers struggle with cryptography apis? In *ICSE* (pp. 935–946).
 - [32] Nair, V., & Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *ICML* (pp. 807–814).
 - [33] Nassif, H., Mohtarami, M., & Glass, J. (2016). Learning semantic relatedness in community question answering using neural models. ACL, (p. 137).

- [34] Nicosia, M., Filice, S., Barrón-Cedeno, A., Saleh, I., Mubarak, H., Gao, W., Nakov, P., Da San Martino, G., Moschitti, A., Darwish, K. et al. (2015). Qcri: Answer selection for community question answeringexperiments for arabic and english. In *SemEval* (pp. 203–209). volume 15.
- [35] Niu, H., Keivanloo, I., & Zou, Y. (2016). Learning to rank code examples for code search engines. *Empirical Software Engineering*, (pp. 1–33).
 - [36] Palomera, D., & Figueroa, A. (2017). Leveraging linguistic traits and semisupervised learning to single out informational content across how-to community question-answering archives. *Information Sciences*, 381, 20–32.
- [37] Petrosyan, G., Robillard, M. P., & De Mori, R. (2015). Discovering information explaining api types using text classification. In *ICSE* (pp. 869–879).
 - [38] Quoc, C., & Le, V. (2007). Learning to rank with nonsmooth cost functions. *NIPS*, 19, 193–200.
 - [39] Raghothaman, M., Wei, Y., & Hamadi, Y. (2016). Swim: Synthesizing what i mean. In *ICSE*.
 - [40] Raja, M. A. Z., Shah, F. H., Alaidarous, E. S., & Syam, M. I. (2017). Design of bio-inspired heuristic technique integrated with interior-point algorithm to analyze the dynamics of heartbeat model. *Applied Soft Computing*, 52, 605–629.

[41] Robertson, S., Walker, S., Jones, S., Hancock-Beaulieu, M. M., & Gatford, M.

(1995). Okapi at trec–3. In *Overview of the Third Text REtrieval Conference* (*TREC–3*) (pp. 109–126). Gaithersburg, MD: NIST.

- [42] Robillard, M. P., & Deline, R. (2011). A field study of api learning obstacles. *Empirical Software Engineering*, 16, 703–732.
- [43] Rohani, V. A., Kasirun, Z. M., Kumar, S., & Shamshirband, S. (2014). An effec-

675

690

695

670

tive recommender algorithm for cold-start problem in academic social networks. *Mathematical Problems in Engineering*, 2014.

- [44] Sakai, T., Ishikawa, D., Kando, N., Seki, Y., Kuriyama, K., & Lin, C.-Y. (2011).
 Using graded-relevance metrics for evaluating community qa answer selection. In *WSDM* (pp. 187–196).
- [45] Schmidhuber, J. (2015). Deep learning in neural networks: An overview. Neural Networks, 61, 85–117.
 - [46] Severyn, A., & Moschitti, A. (2015). Learning to rank short text pairs with convolutional deep neural networks. In *SIGIR* (pp. 373–382).
- [47] Severyn, A., & Moschitti, A. (2016). Modeling relational information in
 question-answer pairs with convolutional neural networks. *arXiv preprint arXiv:1604.01178*, .
 - [48] Singh, P., & Simperl, E. (2016). Using semantics to search answers for unanswered questions in q&a forums. In WWW (pp. 699–706).
 - [49] Sun, R., Cui, H., Li, K., Kan, M.-Y., & Chua, T.-S. (2005). Dependency relation matching for answer selection. In *SIGIR* (pp. 651–652).
 - [50] Sun, Y., Chen, Y., Wang, X., & Tang, X. (2014). Deep learning face representation by joint identification-verification. In *Advances in neural information processing systems* (pp. 1988–1996).
 - [51] Surdeanu, M., Ciaramita, M., & Zaragoza, H. (2008). Learning to rank answers on large online qa collections. In ACL (pp. 719–727). volume 8.

- [52] Tan, M., dos Santos, C., Xiang, B., & Zhou, B. (2016). Improved representation learning for question answer matching. In ACL (pp. 464–473).
- [53] Toba, H., Ming, Z.-Y., Adriani, M., & Chua, T.-S. (2014). Discovering high quality answers in community question answering archives using a hierarchy of classifiers. *Information Sciences*, 261, 101–115.

705

- [54] Turian, J., Ratinov, L., & Bengio, Y. (2010). Word representations: a simple and general method for semi-supervised learning. In ACL (pp. 384–394).
- [55] Van Nguyen, T., Nguyen, A. T., & Nguyen, T. N. (2016). Characterizing api elements in software documentation with vector representation. In *ICSE* (pp. 749– 751).
- [56] Wu, Q., Burges, C. J., Svore, K. M., & Gao, J. (2010). Adapting boosting for information retrieval measures. *Information Retrieval*, 13, 254–270.
- [57] Xue, X., Jeon, J., & Croft, W. B. (2008). Retrieval models for question and answer archives. In *SIGIR* (pp. 475–482).
- ⁷¹⁰ [58] Yan, R., Song, Y., & Wu, H. (2016). Learning to respond with deep neural networks for retrieval-based human-computer conversation system. In *SIGIR* (pp. 55–64).
 - [59] Yang, C., Bai, L., Zhang, C., Yuan, Q., & Han, J. (2017). Bridging collaborative filtering and semi-supervised learning: A neural approach for poi recommendation. In *SIGKDD* (pp. 1245–1254).
 - [60] Yao, Y., Tong, H., Xie, T., Akoglu, L., Xu, F., & Lu, J. (2015). Detecting highquality posts in community question answering sites. *Information Sciences*, 302, 70–82.
 - [61] Yen, S.-J., Wu, Y.-C., Yang, J.-C., Lee, Y.-S., Lee, C.-J., & Liu, J.-J. (2013).
- A support vector machine-based context-ranking model for question answering.
 Information Sciences, 224, 77–87.

- [62] Yu, L., Hermann, K. M., Blunsom, P., & Pulman, S. (2014). Deep learning for answer sentence selection. arXiv preprint arXiv:1412.1632, .
- [63] Zhou, G., Cai, L., Zhao, J., & Liu, K. (2011). Phrase-based translation model for question retrieval in community question answer archives. In ACL (pp. 653–662).
- [64] Zhou, G., He, T., Zhao, J., & Hu, P. (2015). Learning continuous word embedding with metadata for question retrieval in community question answering. In ACL (pp. 250–259).
- [65] Zhou, G., Liu, Y., Liu, F., Zeng, D., & Zhao, J. (2013). Improving question re-

725

- trieval in community question answering using world knowledge. In *IJCAI* (pp. 2239–2245). volume 13.
- [66] Zhou, G., Zhou, Y., He, T., & Wu, W. (2016). Learning semantic representation with neural networks for community question answering retrieval. *Knowledge-Based Systems*, 93, 75–83.
- [67] Zou, Y., Ye, T., Lu, Y., Mylopoulos, J., & Zhang, L. (2015). Learning to rank for question-oriented software text retrieval (t). In ASE (pp. 1–11).

@article{L2ALiSX18,

```
author
              = {Jing Li and Aixin Sun and Zhenchang Xing},
   title
              = {Learning to answer programming questions with software documentation
   through social context embedding},
740
   journal
              = {Information Sciences},
              = \{448 - 449\},\
   volume
              = \{36 - 52\},\
   pages
   vear
              = {2018},
              = {https://doi.org/10.1016/j.ins.2018.03.014},
   url
745
              = {10.1016/j.ins.2018.03.014},
   doi
```

}