Searching by Code: a New SearchBySnippet Dataset and SnippeR Retrieval Model for Searching by Code Snippets

Ivan Sedykh, Dmitry Abulkhanov, Nikita Sorokin,

Sergey Nikolenko, and Valentin Malykh

### **Key Contributions**

- novel "search by snippet" code search task
- novel dataset
- novel model

#### **Task Description**

- Classic Code Search:
  - docstring -> function body
- Clone Detection:
  - function body -> function body
- Search by Snippet:
  - code snippet -> textual description

```
import math
def pythagorus():
    """
    This module takes an input of a, then b, to calculate c in the equation
    sqrt(a^2+b^2) = c
    """
    a,b = int(input("a?")),int(input("b?"))
    c = math.sqrt(a**2+b**2)
    print(c)
```

### Data

- StackOverflow Python-related questions
  - Answers 0
  - Comments 0
  - Other meta-data Ο
- **2 000 000** question-answer pairs
- Test set
  - full-duplicate questions Ο

#### Create ArrayList from array Asked 13 years, 8 months ago Modified 9 days ago Viewed 1.7m times I have an array that is initialized like: 3960 Element[] array = {new Element(1), new Element(2), new Element(3)};



ArrayList<Element> arraylist = ???;





A

0

Share Follow

edited Feb 18, 2020 at 1:53 Unmitigated 46.1k • 7 • 44 • 60 answered Oct 1, 2008 at 14:39 Tom 54.5k • 3 • 26 • 35

#### Train set

 code + traceback is treated like query to find relevant StackOverflow post



"post_id": <u>1767934</u> . post title	
"text": "Why am I getting this error in python ? (httplib) postbody	
<pre><code>if theurl.startswith("http://"): theurl = theurl[7:]</code></pre>	٦
<pre>head = theurl[:theurl.find('/')]</pre>	
<pre>tail = theurl[theurl.find('/'):]</pre>	
response_code = 0	
import httplib	
<pre>conn = httplib.HTTPConnection(head)</pre>	
<pre>conn.request("HEAD",tail)</pre>	
<pre>res = conn.getresponse()</pre>	
<pre>response_code = int(res.status)</pre>	
http://www.garageband.com/mp3cat/.UZCKbS6N4qk/01_Saraenglish.mp3	
Traceback (most recent call last):	
<pre>File "check_data_404.py", line 51, in &lt;module&gt;     run()</pre>	
File "check_data_404.py", line 35, in run	
res = conn.getresponse()	
response.begin()	
File "/usr/lib/python2.6/httplib.py", line 390, in begin version, status, reason = self, read status()	
File "/usr/lib/python2.6/httplib.py", line 354, in _read_status	
httplib RedStatusLine	
d conerd hier	
Does anyone know what "Bad Status Line" is?	
Edit: I tried this for many servers, and many URL's	
and I still get this error?	
The Python Standard Library:	
<a <="" href="https://docs.python.org/2/library/httplib.html" td=""><td></td></a>	
rel="nofollow noreferrer">httplib (Python 2)	
<pre>(called <a <="" href="https://docs.python.org/3/library/http.client.html&lt;/pre&gt;&lt;/td&gt;&lt;td&gt;." td=""></a></pre>	
rel="nofollow noreferrer">http.client in Python 3):	J
best answer	
<pre><blockquote></blockquote></pre>	
<pre><pre><pre><code>exception nttplib.BadStatusLine</code><pre>code&gt;<pre>code&gt;<pre></pre></pre></pre></pre></pre></pre>	
A SUBCLASS OF HIPEXCEPTION. Raised if a server responds with	
<pre>//hlockquotes"</pre>	
(A proprieto cos	

#### SearchBySnippet Dataset

- 900k questions in train set
- 500 questions in test set

```
"text": "How to find out what week number is
current year on June 16th (wk24) with Python? \n",
"score": None,
"probability": None,
"question": None,
"meta": {
    "post_id": "2600775",
    "accepted_answer_id": "0",
    "answer_count": "14",
    "closed_date": None,
    "comment_count": "3",
    "community_owned_date": None,
    "creation_date": "2010-04-08 14:35:53",
    "favorite_count": "56",
    "last_ activity_ date": "2020-03-18 06:55:44",
    "last_edit_ date": "2017-12-13 07:05:52",
    "last_editor_display name": None,
    "last_editor_user_id": "6790377",
    "owner_user_id": "311996",
    "parent_id": "0",
    "post_type_id": "1",
    "score": "334",
    "tags": "<python><datetime><week-number>",
    "title": "How to get week number in Python?",
    "view_count": "265693",
    "python_tag": True
5,
"embedding": None,
"id": "13ca360cb9c37d421e70e0d09c9a3e7dfil"
```

Figure 2: Data sample from the public dump.

#### Model

- Single Encoder Retriever
- Based on PLM
  - We tried several, the best is GraphCodeBERT
- □ Hard-negative mining for self-training

$$\mathcal{L}(q, d^+, D^-) = -\log \frac{e^{\operatorname{score}(q, d^+)}}{\sum\limits_{j=1}^n e^{\operatorname{score}(q, d^-_j)} + e^{\operatorname{score}(q, d^+)}}$$



Figure 1: Overview of the problem setting and system design.

# Negative Hard-Negative





#### Self-training

- the model is used to obtain hard-negatives for the next training iteration



Figure 2: Self-training framework.

#### Results

- existing models lose to simple BM25
- SnippeR outperforms BM25

	Recall			
Model	@5	@10	@20	@50
GraphCodeBERT (Guo et al., 2021)	0.001	0.001	0.002	0.009
CodeBERT (Feng et al., 2020)	0.001	0.006	0.010	0.013
SynCoBERT (Wang et al., 2021)	0.006	0.010	0.013	0.020
GraphCodeBERT (+CSN)	0.161	0.221	0.280	0.367
BM25 (Robertson et al., 1994)	0.311	0.406	0.474	0.562
SnippeR	0.338	0.451	0.536	0.657

Table 3: Results on SearchBySnippet.

#### Conclusion

- Novel task is suggested to community.
- Novel significantly different dataset.
- Novel model effective for this task.

P.S. BM25 is still the best first choice!



## **Problem Description**

- □ Information retrieval
  - **Query**: code + traceback
  - Documents: answers from StackOverflow



## Model

- □ Single Encoder Retriever
- Based on PLM
  - We tried several, the best is GraphCodeBERT
- □ Hard-negative mining for self-training



## GraphCodeBERT



## GraphCodeBERT



## Model

- □ Single Encoder Retriever
- Based on PLM
  - We tried several, the best is GraphCodeBERT
- □ Hard-negative mining for self-training



### Self-Training



### **Achieved Results**

- □ We achieved 5% improvement over BM25
- □ We achieved more than 50% of Recall@10
- □ We transferred our model to Cloud IDE product

Model / Metric	Recall@5	Recall@10	Recall@20	Recall@50
GraphCodeBERT [12]	0.001	0.001	0.002	0.009
CodeBERT [8]	0.001	0.006	0.010	0.013
SynCoBERT [28]	0.006	0.010	0.013	0.020
GraphCodeBERT (CodeSearchNet)	0.161	0.221	0.280	0.367
BM25 [27]	0.311	0.406	0.474	0.562
CodeSentri (ours)	0.338	0.451	0.536	0.657

TABLE V

EVALUATION OF CODESENTRI AND THE BASELINE MODELS ON SEARCHBYSNIPPET DATASET.