# APIzation: Generating Reusable APIs from StackOverflow Code Snippets

[ASE' 21] Valerio Terragni etc.

Lingyu Zhang

2023/03/21

# Outline

- Motivation

- Understanding Real Word APIzations

- APIZATOR: An Automated APIzation Tool for Java Code Snippets

- Evaluation

# Motivation

- Code snippet (CS) from StackOverflow is often incomplete for immediate reuse
  - Lack of type declarations
  - Lack of a well-formed method declaration

```java
 1  // We know week number and year.
 2  int week = 3;
 3  int year = 2010;
 4
 5  // Get calendar, clear it and set week number and year.
 6  Calendar calendar = Calendar.getInstance();
 7  calendar.clear();
 8  calendar.set(Calendar.WEEK_OF_YEAR, week);
 9  calendar.set(Calendar.YEAR, year);
10
11  // Now get the first day of week.
12  Date date = calendar.getTime();
```

# APIzation

- Activity of creating an API for CSs without a well-formed method declaration

```
 1  // We know week number and year.
 2  int week = 3;
 3  int year = 2010;
 4
 5  // Get calendar, clear it and set week number and year.
 6  Calendar calendar = Calendar.getInstance();
 7  calendar.clear();
 8  calendar.set(Calendar.WEEK_OF_YEAR, week);
 9  calendar.set(Calendar.YEAR, year);
10
11  // Now get the first day of week.
12  Date date = calendar.getTime();
```

⬇ Automate?

```
 1  import java.util.Calendar;
 2  import java.util.Date;
 3
 4  public class Human2109186 {
 5      public static Date getFirstDayOfWeek(int week, int year) {
 6          Calendar calendar = Calendar.getInstance();
 7          calendar.clear();
 8          calendar.set(Calendar.WEEK_OF_YEAR, week);
 9          calendar.set(Calendar.YEAR, year);
10          return calendar.getTime();
11      }
12  }
```

# Process of APIzation

1. Choose modifiers and a method name

```java
public static _ getFirstDayOfWeek(_){
    // We know week number and year.
    int week = 3;
    int year = 2010;

    // Get calendar, clear it and set week number and year.
    Calendar calendar = Calendar.getInstance();
    calendar.clear();
    calendar.set(Calendar.WEEK_OF_YEAR, week);
    calendar.set(Calendar.YEAR, year);

    // Now get the first day of week.
    Date date = calendar.getTime();
}
```

# Process of APIzation

2. Recover missing declarations

```java
import java.util.Calendar;
import java.util.Date;

public static _ getFirstDayOfWeek(_){
    // We know week number and year.
    int week = 3;
    int year = 2010;

    // Get calendar, clear it and set week number and year.
    Calendar calendar = Calendar.getInstance();
    calendar.clear();
    calendar.set(Calendar.WEEK_OF_YEAR, week);
    calendar.set(Calendar.YEAR, year);

    // Now get the first day of week.
    Date date = calendar.getTime();
}
```

# Process of APIzation

3. Extract intended input parameters

```java
import java.util.Calendar;
import java.util.Date;

public static _ getFirstDayOfWeek(int week, int year){
    // Get calendar, clear it and set week number and year.
    Calendar calendar = Calendar.getInstance();
    calendar.clear();
    calendar.set(Calendar.WEEK_OF_YEAR, week);
    calendar.set(Calendar.YEAR, year);

    // Now get the first day of week.
    Date date = calendar.getTime();
}
```

# Process of APIzation

4. Extract output

```java
import java.util.Calendar;
import java.util.Date;

public static Date getFirstDayOfWeek(int week, int year){
    // Get calendar, clear it and set week number and year.
    Calendar calendar = Calendar.getInstance();
    calendar.clear();
    calendar.set(Calendar.WEEK_OF_YEAR, week);
    calendar.set(Calendar.YEAR, year);
    return calendar.getTime();
}
```

# Understanding Real Word APIzations

- Data collection approach
    - Explicit StackOverflow link
    - Type 3 code clone
    - Manual check


- 135 $< CS, API >$ pairs reference 509 variables

# Findings on method parameters

- PATT-notdecl



- PATT-const

# Findings on return statements

- ## PATT-latest



- ## PATT-syso

# Choose Modifiers and A Method Name

- Modifiers: `public static`
  - API must be accessible by any other class
  - Avoiding instantiating objects for invoking the API

- Method name: Part-of-Speech (POS) Tagger
  - Generating from the title of the corresponding StackOverflow page
    - Assigning parts of speech to each word in the title
    - Combining the main "verb" and the corresponding "direct object"

How to get first day of a given week number in Java

Asked 13 years, 1 month ago    Modified 1 year, 6 months ago    Viewed 50k times

➡️    getDay

# Recover Missing Declarations

- Type Declarations: CSNIPPEX*
  - A greedy algorithm based on the clustering hypothesis

| File | IOException | PrintWriter | Document | Jsoup |
|------|-------------|-------------|----------|-------|
| java.io | java.io | java.io | org.bson ❌ | org.jsoup |
| scala.. | com.sun.. | | org.jdom2 | |
| org.specs.. | net.kuujo.. | | org.jsoup.nodes | |
| .... | ... | | .... | |

*: Valerio Terragni, Yepang Liu, Shing-Chi Cheung,CSNIPPEX: automated synthesis of compilable code snippets from Q&A sites.
  ISSTA 2016: 118-129

2023/6/12

13

# Recover Missing Declarations

- Variable Declaration: BAKER*
  - Identifying the most plausible type of $v$ by leveraging the usages of $v$ in the API



The answer above is almost 100% correct. It will fail with unicode.

```
 5    1 MessageDigest digest;
      2 try {
      3     digest = MessageDigest.getInstance("MD5");
      4     byte utf8_bytes[] = tag_xml.getBytes();
      5     digest.update(utf8_bytes,0,utf8_bytes.length);
      6     hash = new BigInteger(1, digest.digest()).toString(16);
      7 }
      8 catch (NoSuchAlgorithmException e) {
      9     e.printStackTrace();
     10 }
```

Need the length from the byte array not the string.

share improve this answer

answered

$\longrightarrow$ String tax_xml

*: Siddharth Subramanian, Laura Inozemtseva, Reid Holmes. Live API documentation. ICSE 2014: 643-652

# Extract intended input parameters

- PATT-notdecl
  - Undeclared variables are input parameters

$$
\begin{aligned}
&/\ast \qquad\qquad\qquad \textbf{PATT-notdecl} \qquad\qquad\qquad\qquad \ast/\\
&\textbf{else if } errors \subseteq \texttt{missing-variable-decl} \textbf{ then}\\
&\quad \textbf{for } v \in (errors \cap \texttt{missing-variable-decl}) \textbf{ do}\\
&\qquad \langle \tau, imports, classpath \rangle \leftarrow \textsc{RecoverVarType}(v, API,\\
&\qquad JARs, imports, classpath)\\
&\qquad \mathcal{T}[v] \leftarrow \tau\\
&\qquad \text{add } \langle \tau, v \rangle \text{ to } API.\textit{parameter-list}
\end{aligned}
$$

# Extract intended input parameters

- PATT-const
  - Variables with const value are input parameters

```
/*                          PATT-const                          */
LP-VARS ← GetLoopChangingVars(API.method-body)
for s_i ∈ API.method-body do
    case s_i : τ v = ε do      // Variable decl. and init.
        T[v] ← τ
        add v to ALREADY-INIT-VARS
        if IsHardCoded(τ, ε) ∧ v ⊄ LP-VARS then
            add ⟨τ, v⟩ to API.parameter-list
            remove s_i from API.method-body

    case s_i : τ v do                // Variable declaration
        ⟨T[v], S[v]⟩ ← ⟨τ, s_i⟩

    case s_i : v = ε do              // Variable assignment
        if v ∉ ALREADY-INIT-VARS then
            add v to ALREADY-INIT-VARS
            if IsHardCoded(τ, ε) ∧ v ∉ LP-VARS then
                add ⟨T[τ], v⟩ to API.parameter-list
                remove s_i from API.method-body
                remove S[v] from API.method-body
```

# Extract output

- PATT-latest

```
/*                          PATT-latest                          */
case sₙ : τ v = ε do          // Variable decl. and init.
    API.return-type ← τ
    replace sₙ in API.method-body with return ε;

case sₙ : v = ε do                  // Variable assignment
    API.return-type ← T[v]
    replace sₙ in API.method-body with return ε;
```

The monospace code block contains mathematical notation. Rendered precisely:

$$/* \qquad \textbf{PATT-latest} \qquad */$$

**case** $s_n : \tau\ v = \epsilon$ **do**    // Variable decl. and init.
  $API.return\text{-}type \leftarrow \tau$
  replace $s_n$ in *API.method-body* with return $\epsilon$;

**case** $s_n : v = \epsilon$ **do**    // Variable assignment
  $API.return\text{-}type \leftarrow \mathcal{T}[v]$
  replace $s_n$ in *API.method-body* with return $\epsilon$;

- PATT-syso

$$/* \qquad \textbf{PATT-syso} \qquad */$$

**case** $s_n$ : `System.out.println(string-literal + ` $\epsilon$ `)` $\lor$
  `System.out.println(` $\epsilon$ `)` **do**
  $API.return\text{-}type \leftarrow \textsc{GetTypeOfExp}(\epsilon, imports, classpath)$
  replace $s_n$ in *API.method-body* with return $\epsilon$;

**otherwise do**
  $API.return\text{-}type \leftarrow$ `void`
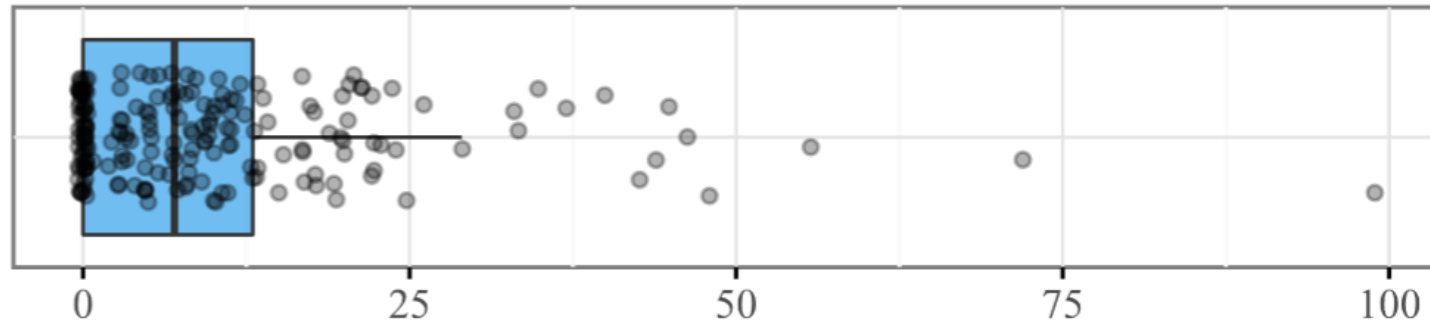
# Evaluation

- RQ1: Identical APIs



Fig. 3.  Distribution of the number of AST differences.

**RQ1 – In summary:**  APIZATOR generated 63 (31.50 %) APIs identical (including the method-body and import declarations) to the human-produced ones.

# Evaluation

- RQ2: Method Parameters

## TABLE I
### RQ2 ANALYSIS AND COMPARISON OF THE HUMAN- ($P_H$) AND APIZATOR-PRODUCED ($P_A$) PARAMETER LISTS

| Param. $\lvert P_H \rvert$ | Human APIs | $P_H \equiv P_A$ Count | % | $\lvert P_H \setminus P_A \rvert$ Mean | Min | Mdn | Max | $\lvert P_H \cap P_A \rvert$ Mean | Min | Mdn | Max | $\lvert P_A \setminus P_H \rvert$ Mean | Min | Mdn | Max | Jaccard Distance (JD) Mean | Min | Mdn | Max |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 58 | 45 | 77.59 | – | – | – | – | – | – | – | – | 0.36 | 0.00 | 0.00 | 5.00 | 0.22 | 0.00 | 0.00 | 1.00 |
| 1 | 93 | 60 | 64.52 | 0.32 | 0.00 | 0.00 | 1.00 | 0.68 | 0.00 | 1.00 | 1.00 | 0.13 | 0.00 | 0.00 | 2.00 | 0.34 | 0.00 | 0.00 | 1.00 |
| 2 | 35 | 7 | 20.00 | 1.14 | 0.00 | 1.00 | 2.00 | 0.86 | 0.00 | 1.00 | 2.00 | 0.29 | 0.00 | 0.00 | 2.00 | 0.58 | 0.00 | 0.50 | 1.00 |
| $\geq 3$ | 14 | 1 | 7.14 | 2.86 | 0.00 | 3.00 | 6.00 | 0.64 | 0.00 | 0.00 | 4.00 | 0.21 | 0.00 | 0.00 | 1.00 | 0.82 | 0.00 | 1.00 | 1.00 |
| Total ($\geq 0$) | 200 | 113 | 56.50 | 0.77 | 0.00 | 0.50 | 6.00 | 0.72 | 0.00 | 1.00 | 4.00 | 0.23 | 0.00 | 0.00 | 5.00 | 0.38 | 0.00 | 0.00 | 1.00 |

**RQ2 – In summary:** APIZATOR generated 113 (56.50 %) APIs with identical parameter lists to the human-produced ones.

# Evaluation

- RQ3: Return Statements

**TABLE II**
**RQ3 RETURN STATEMENTS COMPARISON**

| $API_H$ | $API_A$ | Count | % | Count | % |
|---|---|---|---|---|---|
| | **Return Type** | | | **Equivalent Return Type and Statements** | |
| void | void | 63 | 31.50 | 63 | 100.00 |
| void | not void | 2 | 1.00 | – | – |
| not void | void | 72 | 36.00 | – | – |
| not void | not void | 63 | 31.50 | 52 | 82.54 |
| Total | | 200 | | 115 | |

**RQ3 – In summary:** APIZATOR generated 115 (57.50 %) APIs with identical return statements to the human-produced ones.

# Thanks

Comments are welcome!