

### **String Test Data Generation for Java Programs**

Miaomiao Wang, Baoquan Cui, Jiwei Yan, Jun Yan, Jian Zhang

Technology Center of Software Engineering, Institute of Software, Chinese Academy of Sciences, Beijing, China State Key Laboratory of Computer Science, Institute of Software, Chinese Academy of Sciences, Beijing, China University of Chinese Academy of Sciences, Beijing, China

**ISSRE 2022** 



### Overview

### Introduction

- String Test Data Generation
- Existing Work

### Preliminary

- String API Usage
- Motivating Example

#### Our Approach

- API-Regex Mapping
- String API Invocation Sequence Extraction
- Regex Generation
- $\succ$  Evaluation
  - DataSets and Experimental Results
  - Case Study
- Conclusion

### **1.1 String Test Data Generation**



3

- String operations in Java programs are mostly performed by calling string APIs.
- We need to understand the semantics of the API to generate suitable test data.

```
1 // str is a method parameter
2 String[] splitArray = str.split("#");
3 String splitStr = splitArray[2];
4 String subStr = splitStr.substring(3, 8);
5 if (subStr.equals("class")) {
6 ... //True branch: omitted
7 } else { ... //False branch: omitted }
```

complex semantics for single API

multiple string API combination

### **1.2 Existing Work**

#### SCAS

#### • **SMT Solvers**([1]-[8]):

- SMT-LIB[9] cannot support the split(...) method and the lastIndexOf(...) method
- The automatic conversion from the string API combinations to SMT formats is not supported.

#### test case generation

- EvoSuite[10]
  - Seeding Strategy [11] : It supports single and simple string APIs to seed data and mutate
- Randoop[12]
  - feedback-directed random test generation
  - It does not intentionally generate string-related data.
- [1] Norn: An SMT Solver for String Constraints
- [2] String constraints with concatenation and transducers solved efficiently
- [3] Towards Constraint Logic Programming over Strings for Test Data Generation
- [4] On Solving Word Equations Using SAT.
- [5] Z3str4: A Solver for Theories over Strings
- [6] cvc5: A Versatile and Industrial-Strength SMT Solver
- [7] A decision procedure for path feasibility of string manipulating programs with integer data type
- [8] Solving string constraints with regex-dependent functions through transducers with priorities and variables
- [9] https://smtlib.cs.uiowa.edu/theories-UnicodeStrings.shtml
- [10] Evosuite: automatic test suite generation for object-oriented software
- [11] Seeding strategies in search-based unit test generation
- [12] Randoop: feedback-directed random testing for java



### Overview

### Introduction

- String Test Data Generation
- Existing Work

### > Preliminary

- String API Usage
- Motivating Example
- Our Approach
  - API-Regex Mapping
  - String API Invocation Sequence Extraction
  - Regex Generation
- $\succ$  Evaluation
  - DataSets and Experimental Results
  - Case Study
- Conclusion

### 2.1 String API Usage



		String	API USAGE			
ΔΡΙ	NO	In	Path	In Branch		
AII	10.	ct.(Total)	ct.(Dedup)	ct.(Total)	ct.(Dedup)	
Single	1	11,536	29	8,304	25	
Combination	2	873	84	753	75	
	3	159	58	132	49	
	4	32	7	27	6	
	5	4	2	3	1	
Sum		12,604	180	9,219	156	

TABLE II String API Usage

Single: a single string API invocation.

Combination: combination of multiple API invocations.

NO: the number of API combinations, a single API involves only one API.

ct. (Total): the usage the API combination is used.

ct. (Dedup): the count the API combination is used after deduplication.

- The combinations of string APIs make up 8.5% (1068/12604).
- 73.1% (9219/12604) of string API invocations or their combinations have affected the branches in the programs.

### **2.2 Motivating Example**





- SMT solvers, EvoSuite and Randoop can not generate suitable test data to cover the true branch.
- The purpose of this paper is to generate regular expressions:
  - cover the branches
  - trigger the potential exception proactively



### Overview

### Introduction

- String Test Data Generation
- Existing Work

### Preliminary

- String API Usage
- Motivating Example

### > Our Approach

- API-Regex Mapping
- String API Invocation Sequence Extraction
- Regex Generation
- $\succ$  Evaluation
  - DataSets and Experimental Results
  - Case Study
- Conclusion



### **Our Approach**



Fig. 1. Overview of Our Approach.

### **3.1 SAIS Extraction**

#### SCAS

#### Algorithm 1: Slice

```
Input: StrPP, condExpr
   Output: snippet
1 containsStringParameter = False;
2 snippet = \{\};
3 variableSet = \{\};
4 stmt = condExpr;
5 snippet.insertAtFirst(stmt);
6 variableSet.addAll(getVariables(stmt));
7 while stmt != NULL do
       for each v \in variableSet do
           if stmt.isAssignStatementOf(v) then
 9
               variableSet.remove(v);
10
               snippet.insertAtFirst(stmt);
11
               variables = getVariables(stmt);
12
               variableSet.addAll(variables);
13
               if isStringParaLocalization(strnt) then
14
                   containsStringParameter = true;
15
               break;
16
       if variableSet.isEmpty() then
17
           break;
18
       stmt = StrPP.getPredecessorOf(stmt);
19
20 if containsStringParameter == False then
       snippet = NULL;
21
22 return snippet
```

#### **Definition III-A1. (StringParaPath)**

A *StringParaPath* is the path in a method which contains at least one string API invocation statement on the local variable of the string parameter.







**Definition III-A2. (String API Invocation Sequence (SAIS))** A **SAIS** is a 2-tuple

**SAIS = <CondExpr, Trajectory>** where the *CondExpr* is a conditional expression in the StringParaPath, and the Trajectory is a k-length list of ordered pairs

 $<(\text{stmt}_1, \text{map}_1), (\text{stmt}_2, \text{map}_2), ..., (\text{stmt}_k, \text{map}_k) >$ where the *stmt*<sub>*i*</sub> in each pair is a string API invocation statement  $\frac{1}{3}$ related to the conditional expression and the  $map_i$  maps the 5 variables in the  $stmt_i$  to their values. 6

### **3.1 SAIS Extraction**



```
// str is a method parameter
String[] splitArray = str.split("#");
String splitStr = splitArray[2];
String subStr = splitStr.substring(3, 8);
if (subStr.equals("class")) {
  ... //True branch: omitted
} else { ... //False branch: omitted }
```

SAIS	CondExpr:	eq==True	
	Trajectory:	$splitStr = str.split(v_4)[i_1]$	$\{v_4:"\#", i_1:2\}$
		substr = splitStr.substring(v)	$(v_2,v_3)\{v_2:3,v_3:8\}$
		$eq = subStr.equals(v_1),$	$\{v_1:$ "class" $\}$
1 1 1		All and a second second	
		and a second	

7



### **3.1 SAIS Extraction**





### **3.2 API-Regex Mapping**



#### • Terminated API.

• The API whose return type is char, int, boolean, char[] or byte[], that is, primitive type or its array type, is considered as a terminated API (Total: 15). In general, its return value is used directly in the conditional expression

#### Non-terminated API.

• The API whose return type is String, CharSequence or String[] is considered as a non-terminated API (Total: 33). Their return values cannot appear directly in conditional expressions.



#### SCAS

14

## **3.2 API-Regex Mapping**

### Definition (RegexWrapper)

A *RegexWrapper* is a 4-tuple

- R is the regular expression according to the string API and the expression on its return value;
- L is the length of the string that the R allows to merge and -1 indicates that there is no limit for the length;
- S is the a special regular expression with which to replace during inference if the any character regular expression is in the suffix of the R;
- $L_{min}$  is the minimum length of a string corresponding to the R;

### Definition (API-Regex Mapping)

A API-Regex Mapping is a 2-tuple

 $\mathcal{M} = \langle APIPair, RegexWrapper \rangle$ where *APIPair* is a pair including the string API invocation and its related conditional expression, and *RegexWrapper* is defined before.

### **3.2 API-Regex Mapping**



#### TABLE IV Partial API-Regex Mapping

APIPair	RegexWrapper				
API Invo. (caller: <i>str</i>   <i>arr</i> )	CondExpr	R	L	S	$L_{min}$
r = str.startswith(v)	r == T	$v[\langle s \rangle S]^*$	-1	NULL	Len(v)
r = str.length()	r = l	$[\S\S]{l}$	-1	NULL	l
$r = str.substring(i_1, i_2)$	NULL	$[\S\]\{i_1\}$	$i_2$ - $i_1$	NULL	$i_1$
r = str.split(v)[i]	NULL	$[v]{i}$	-1	$[^v]$	i
r=arr.length	r > i	$[\s\S]{i+1,}$	-1	NULL	<i>i</i> +1

https://github.com/suoyi123wang/JustinStr

### **3.3 Regex Generation**





### **3.3 Regex Generation**



		D	1 5				
Tool	E-prone Regular Expression						
1001	$[\s\S]{1,4}$	$[^a]{1,3}$	(?i)AbC	$[\langle Q^* \rangle E]{1,3}$			
xeger [4]	×	<b>~</b>	×	×			
Generex [39]	<ul> <li>✓</li> </ul>	<b>~</b>	×	Е			
MutRex [41]	~	<b>v</b>	X	×			
bfgex [22]	×	X	Е	×			
RgxGen [19]	<ul> <li>✓</li> </ul>	<b>~</b>	X	<b>~</b>			
random-string [9]	E	E	×	E			

 TABLE I

 COMPARISON ON JAVA TOOLS GENERATING STRING DATA

 $\checkmark$ : the generated string matches the regex;  $\checkmark$ : not matching; E: an exception occurs during generation.

		- 17
<b>Cover the true branch (line 6)</b>	7	<pre>} else { //False branch: omitted }</pre>
	6	//True branch: omitted
g	5	<pre>if (subStr.equals("class")) {</pre>
RgxGen: "##abcclass"	4	<pre>String subStr = splitStr.substring(3, 8);</pre>
	3	<pre>String splitStr = splitArray[2];</pre>
	2	<pre>String[] splitArray = str.split("#");</pre>
[#]\7\[^#]\3\c]ass	1	<pre>// str is a method parameter</pre>



### Overview

### Introduction

- String Test Data Generation
- Existing Work

### Preliminary

- String API Usage
- Motivating Example

#### Our Approach

- API-Regex Mapping
- String API Invocation Sequence Extraction
- Regex Generation

### Evaluation

- DataSets and Experimental Results
- Case Study

Conclusion

### 4.1 DataSets



#### TABLE V STATISTICS OF THREE DATA SETS

Program	#Cls	#LOC	$M_1$	$M_2$	#RE
23 Assginments	406	9667	45	42	98
289 Solutions	289	4529	217	183	677
closure-compiler-20220202	8,162	218,939	194	153	475
commons-lang3-3.12.0	505	9,600	73	39	211
commons-io-2.11.0	240	5,544	21	17	81
mysql-connector-java-8.0.29	1,345	42,726	56	44	163
poi-ooxml-5.2.2	1,529	47,176	45	31	84
OpenJDK-8u292	25,574	697,752	1,513	1,250	6,697





### **4.2 Experimental Results**

• RQ1: How effectively does Justin-Str characterize the input string parameters?

Length	Branch	Single (%)	Combination (%)				
Length	Dranen	Single (70)	2	3	4	5	
L=10,20	True	100.00	100.00	100.00	100.00	100.00	
(30,40,50)		100100				100100	
L=10	False	99.77	99.41	99.33	98.09	100.00	
L=20	False	99.55	99.14	99.02	99.11	100.00	
L=30	False	99.35	98.84	98.72	99.35	99.99	
L=40	False	99.14	98.47	98.50	99.49	99.99	
L=50	False	98.94	98.20	98.18	99.65	99.98	

- 85x5 = 425 methods,
  - Each one ends with a branch (T/F)

### **4.2 Experimental Results**



21

• **RQ2: Improvement of the branch coverage** 

Program	Seeds	$M_{Cov100\%}$	$M_{Cov}(\%)$	$Max_{Cov}$	Ave <sub>Cov</sub>
closure-compiler	972	58	+18(19%)	+50%	+21%
commons-lang3	441	10	+10(34%)	+57%	+11%
commons-io	128	4	+5(38%)	+43%	+17%
mysql-connector-java	301	16	+5(18%)	+34%	+11%
poi-ooxml	125	10	+2(10%)	+50%	+28%

- +22%(40/186) of the method
- up to +57%, and +17% on average

### **4.2 Experimental Results**

#### • RQ3: Bug Finding

Program	InstinStr	Common	Δ	Time (s)		
Tiogram	Justinou	Common		JustinStr	EvoSuite	
Assignments	38	27	11	339	1627	
Solutions	72	60	12	14	12,453	
closure-compiler	33	21	12	72	2,394	
commons-io	11	1	10	16	80	
commons-lang3	15	3	12	15	709	
mysql-connector-java	7	3	4	23	306	
poi-ooxml	3	0	3	40	-	
OpenJDK	17	0	17	10,228	-	
Total	196	115	81	-	-	

### TABLE VIII COMPARISON OF BUGS FOUND BY JUSTINSTR AND EVOSUITE

# TABLE IXBUGS CONFIRMED OR FIXED BY THE JDK DEVELOPERSe:ArrayIndexOutOfBoundsException (2)

Type:	ArrayIndexOutOrBoundsException (2)
Issue ID	: I4MWI1 (Commit ID), 8279422
Туре:	StringIndexOutOfBoundsException (14)
Issue ID	: 8278186, 8279128, 8279129, 8279198, 8279218, 8279336,
8279341	, 8279342, 8279362, 8279423, **21212bd18(Commit ID),
8279424	, **411a404a9(Commit ID), **8baba7d11(Commit ID)
Туре:	Infinite Loop (1), <b>Issue ID</b> : 8278993



22

### **4.2 Experimental Results**

#### iscas

#### • RQ3: Bug Finding (Order)

•



Fig. 3. Comparison on Number of Bugs Found by JustinStr and EvoSuite under Test Case Order. Bars represent the number of bugs triggered by the test case at current order and poly-lines represent the number of accumulated bugs triggered by the test cases before and within the order.

The first two test cases in JustinStr can trigger 74% of the bugs, while EvoSuite only triggers 36%.

### 4.3 Case Study

#### SCAS

24



https://bugs.java.com/bugdatabase/view\_bug.do?bug\_id=JDK-8278186 https://bugs.java.com/bugdatabase/view\_bug.do?bug\_id=JDK-8278993



### Overview

### Introduction

- String Test Data Generation
- Existing Work

### Preliminary

- String API Usage
- Motivating Example

#### Our Approach

- API-Regex Mapping
- String API Invocation Sequence Extraction
- Regex Generation
- $\succ$  Evaluation
  - DataSets and Experimental Results
  - Case Study

### Conclusion

### Conclusion



- **API-Regex Mapping**. We build a mapping from 48 string APIs to regular expressions based on their semantics, which are equivalent or approximate.
- **Inference Algorithm**: We design an inference algorithm that can characterize the string parameter with regular expressions under the combination of its string API invocations.
- String Test Data Generation Tool. We developed an automatic tool JustinStr to output the string test data for test cases generation.

Future work:

- Convert the string API with the symbolic value into regular expressions.
- Support for converting more string APIs to regular expressions.



# Q & A