

Grammar-Agnostic Symbolic Execution by Token Symbolization

Zhenbang Chen (<u>zbchen@nudt.edu.cn</u>)

Joint work with Weiyu Pan, Guofeng Zhang, Yunlai Luo, Yufeng Zhang, and Ji Wang



ACM SIGSOFT International Symposium on Software Testing and Analysis



Parsing







Parsing











An input grammar usually exists



Complex Parsing Code

An input grammar usually exists











- Random testing
- Coverage-oriented fuzzing
 - AFL, JQF, ...

- Random testing
- Coverage-oriented fuzzing
 - AFL, JQF, ...

Not efficient

• Random testing

- Coverage-oriented fuzzing
 - AFL, JQF, ...
- Grammar-directed fuzzing
 - Black-box fuzzing[ASE'19], Zest[ISSTA'19], ...

Random testing

- Coverage-oriented fuzzing
 AFL, JQF, ...
- Grammar-directed fuzzing
 - Black-box fuzzing[ASE'19], Zest[ISSTA'19], ...
- Grammar-directed symbolic execution
 - Godefroid et al. [PLDI'08], CESE[ASE'07], ...

execution CESE[ASE'07]....

20

• Random testing

- Coverage-oriented fuzzing • AFL, JQF, ...
- Grammar-directed fuzzing
 - Black-box fuzzing[ASE'19], Zest[ISSTA'19], ...
- Grammar-directed symbolic execution
 - Godefroid et al. [PLDI'08], CESE[ASE'07], ...

Need Input Grammar



• Random testing

- Coverage-oriented fuzzing • AFL, JQF, ...
- Grammar-directed fuzzing
 - Black-box fuzzing[ASE'19], Zest[ISSTA'19], ...
- Grammar-directed symbolic execution
 - Godefroid et al. [PLDI'08], CESE[ASE'07], ...





Symbolic Execution of Complex Parsing Programs without Grammar



Symbolic Execution of Complex Parsing Programs without Grammar



Many inputs are rejected



Symbolic Execution of Complex Parsing Programs without Grammar



Symbolic Execution of Complex Parsing Programs without Grammar



Analyze the application logic code is even harder



Key Insight





Key Insight





Key Insight

Complex Parsing Code : Syntactic Analysis AST Token Sequence





Token provides an abstraction for the input grammar







Key Idea

Symbolize token values





Technique Problems

Technique Problems

Input generation for a token sequence




Technique Problems

Input generation for a token sequence

Analyze the application code in priority





Technique Problems

Input generation for a token sequence



Summarize the tokenization code first



Technique Problems

Input generation for a token sequence



Summarize the tokenization code first

Record byte-level path constraint in application logic





















- $\langle Expr \rangle$::= $\langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$ (ID) ::= 'a' | 'b' | ... | 'y' | 'z'(Number) ::= '00' | '01' | ... | '98' | '99'
 - ⟨Op⟩ ::= '*' | '+'

- ⟨Op⟩ ::= '*' | '+'



 $\langle Expr \rangle$::= $\langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$ (ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99'

II+I2 3I+a 49+z II*a ...

- ⟨Op⟩ ::= '*' | '+'
- - $||+|a 0*|| a-|| a+|| \cdots$

 $\langle Expr \rangle$::= $\langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$ (ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99'

II+I2 3I+a 49+z II*a ...

- (ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99' ⟨Op⟩ ::= '*' | '+'
 - 2 // inputReader's type is Reader 3 inputReader = new StringReader(a); 4 parseExpr(); 5 // application logic starts 6 if (a.charAt(a.length() - 1) == 'z') { assert(false); //bug 7 8 9 }

 $\langle Expr \rangle$::= $\langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$

public void entry(String a) throws ParseException{

(ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99' ⟨Op⟩ ::= '*' | '+'



 $\langle Expr \rangle$::= $\langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$



11	<pre>void parseExpr() throws</pre>
12	<pre>int token = getNextTok</pre>
13	<pre>if (token == T_NUM){</pre>
14	<pre>parse0p();</pre>
15	return;
16	<pre>} else if (token == T_</pre>
17	<pre>if (getNextToken() ==</pre>
18	}
19	throw new ParseExcepti
20	}

ParseException { ();

ID){ T_EOF) return;

on();

32	<pre>int getNextToken() throws ParseException {</pre>
33	<pre>int res = inputReader.read();</pre>
34	<pre>if (res == -1) return T_EOF;</pre>
35	<pre>char c = (char) res;</pre>
36	if (c >= '*' && c <= '+'){
37	<pre>return T_OP;</pre>
38	<pre>} else if (c >= 'a' && c <= 'z'){</pre>
39	<pre>return T_ID;</pre>
40	<pre>} else if (c >= '0' && c <= '9'){</pre>
41	<pre>char next_c = (char) inputReader.read();</pre>
42	<pre>if (next_c >= '0' && next_c <= '9'){</pre>
43	<pre>return T_NUM;</pre>
44	}
45	}
46	<pre>throw new ParseException();</pre>
47	}

$$\langle Expr \rangle$$
 ::=
 $\langle ID \rangle$::=
 $\langle Number \rangle$::=
 $\langle Op \rangle$::=

(ID) | (Number) | (Number) (Op) (Expr) 'a' | 'b' | ... | 'y' | 'z' '00' | '01' | ... | '98' | '99' '*' | '+'

> T_{D} T_NUM T OP



Pure DSE

• Pure DSE

First Path's Path Condition **Input:** "12+13"

```
32 int getNextToken() throws ParseException {
int res = inputReader.read();
34 if (res == -1) return T_EOF;
35 char c = (char) res;
  if (c >= '*' && c <= '+'){
37 return T_OP;
38 } else if (c >= 'a' && c <= 'z'){</pre>
39 return T_ID;
40 } else if (c >= '0' && c <= '9'){
41 char next_c = (char) inputReader.read();
42 if (next_c >= '0' && next_c <= '9'){
  return T_NUM;
44 }
46 throw new ParseException();
47 }
```

 $a[4] \neq 'z'$

 $a[0] \ge '*' \land a[0] > '+' \land a[0] < 'a' \land$ $a[0] \ge '0' \land a[0] \le '9' \land a[1] \ge '0' \land a[1] \le '9' \land a[1]$ $a[2] \ge '*' \land a[2] \le '+' \land$ $a[3] \ge '*' \land a[3] > '+' \land a[3] < 'a' \land$ $a[3] \ge '0' \land a[3] \le '9' \land a[4] \ge '0' \land a[4] \le '9' \land a[4]$





• Pure DSE

Input: "12+13"

 $a[0] \geq '*' \wedge a[0]$ $a[0] \geq '0' \wedge a[0]$ $a[2] \geq '*' \wedge a[2]$ $a[3] \geq '*' \wedge a[3]$ $a[3] \geq '0' \wedge a[3]$ $] \neq 'z'$

First Path's Path Condition

$$\begin{array}{l} 0] > '+' \wedge a[0] < 'a' \wedge \\ 0] \leq '9' \wedge a[1] \geq '0' \wedge a[1] \leq '9' \wedge \\ 2] \leq '+' \wedge \\ 3] > '+' \wedge a[3] < 'a' \wedge \\ 3] \leq '9' \wedge a[4] \geq '0' \wedge a[4] \leq '9' \wedge \end{array}$$

Negation in DFS style

• Pure DSE

Input: "12+13"

- $a[0] \ge '*' \land a[0] >$ $a[0] \geq '0' \wedge a[0] \leq$ $a[2] \geq '*' \wedge a[2] \leq$ $a[3] \ge '*' \land a[3] > '+' \land a[3] < 'a' \land$ $a[3] \ge '0' \land a[3] \le '9' \land a[4] \ge '0' \land a[4] \le '9' \land a[4]$ a[4] = 'z'
- First Path's Path Condition

$$\begin{array}{l} & + & | \wedge a[0] < & | a & | \wedge \\ & \leq & | 9 & | \wedge a[1] \geq & | 0 & | \wedge a[1] \leq & | 9 & | \wedge \\ & \leq & | + & | \wedge \\ & & | + & | \wedge a[3] < & | a & | \wedge \end{array}$$

• Pure DSE

Input: "12+13" First Path

- $\begin{array}{l} a[0] \geq '*' \wedge a[0] > '+' \wedge a[0] < 'a' \wedge \\ a[0] \geq '0' \wedge a[0] \leq '9' \wedge a[1] \geq '0' \wedge a[1] \leq '9' \wedge \\ a[2] \geq '*' \wedge a[2] \leq '+' \wedge \\ a[3] \geq '*' \wedge a[3] > '+' \wedge a[3] < 'a' \wedge \\ a[3] \geq '0' \wedge a[3] \leq '9' \wedge a[4] \geq '0' \wedge a[4] \leq '9' \wedge \\ a[4] = 'z' \end{array}$
- First Path's Path Condition

Unsatisfiable

- $\langle Expr \rangle ::= \langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle$ (ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99' • Pure DSE ⟨Op⟩ ::= '*' | '+'
- Input: "12+13"
 - $a[0] \ge '*' \land a[0] >$ $a[0] \geq '0' \wedge a[0] \leq$ $a[2] \geq '*' \wedge a[2] \leq$ $a[3] \ge '*' \land a[3] > '+' \land a[3] < 'a' \land$
 - $a[3] \ge '0' \land a[3] \le '9' \land a[4] \ge '0' \land a[4] > '9'$

First Path's Path Condition

$$(+) + (+) \wedge a[0] < (a' \wedge a') \\ (9' \wedge a[1] \ge (0' \wedge a[1]) \le (9' \wedge a') \\ (+) \wedge a[1] \ge (0' \wedge a[1]) \le (9' \wedge a')$$

- (Expr) ::= (ID) | (Number) | (Number) (Op) (Expr)
 (ID) ::= 'a' | 'b' | ... | 'y' | 'z'
 (Number) ::= '00' | '01' | ... | '98' | '99'
 (Op) ::= '*' | '+'
- Input: "12+13" First Path
 - $a[0] \ge '*' \land a[0] >$ $a[0] \ge '0' \land a[0] \le$ $a[2] \ge '*' \land a[2] \le$
 - $a[2] \ge * \land a[2] \ge + \land A[3] \ge ' * \land a[3] < 'a' \land A[3] > ' + ' \land a[3] < 'a' \land A[3] < (a' \land A[3]) < ($
 - $a[3] \ge '0' \land a[3] \le '9' \land a[4] \ge '0' \land a[4] > '9'$
- New Input: "12+1z"

First Path's Path Condition

$$|+' \wedge a[0] < |a' \wedge a|$$

$$|9' \wedge a[1] \ge |0' \wedge a[1] \le |9' \wedge a|$$

- (Expr) ::= (ID) | (Number) | (Number) (Op) (Expr)
 (ID) ::= 'a' | 'b' | ... | 'y' | 'z'
 (Number) ::= '00' | '01' | ... | '98' | '99'
 (Op) ::= '*' | '+'
- Input: "12+13" First Path
 - $a[0] \ge '*' \land a[0] \ge$ $a[0] \ge '0' \land a[0] \le$ $a[2] \ge '*' \land a[2] \le$
 - $a[3] \ge '*' \land a[3] >$
 - $a[3] \ge '0' \land a[3] \le$
- New Input: "12+1z"

First Path's Path Condition

$$+ ' \wedge a[0] < 'a' \wedge$$

$$+ ' \wedge a[1] \ge '0' \wedge a[1] \le '9' \wedge$$

$$+ ' \wedge$$

$$+' \wedge a[3] < 'a' \wedge$$

$$9' \wedge a[4] \geq 0' \wedge a[4] > 9'$$

New input will be rejected

• Pure DSE

〈Expr〉	::=	$\langle D \rangle$	{Nu
$\langle D \rangle$::=	'a'	'b'
(Number)	::=	'00'	'01
$\langle Op \rangle$::=	'*'	'+'

Six iterations to cover all the parsing code

• Pure DSE

public void entry(String a) throws ParseException{ // inputReader's type is Reader 2 inputReader = new StringReader(a); 3 parseExpr(); 4 // application logic starts 5 if (a.charAt(a.length() - 1) == 'z') { assert(false); //bug 9

```
\langle Expr \rangle ::= \langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle
                            (ID) ::= 'a' | 'b' | ... | 'y' | 'z'
```

Impossible to trigger the bug with 5-length strings



Grammar-Agnostic DSE

• Grammar-agnostic DSE

```
int getNextToken() throws ParseException {
   int res = inputReader.read();
  if (res == -1) return T_EOF;
34
   char c = (char) res;
35
   if (c >= '*' && c <= '+'){
36
    return T_OP;
37
   } else if (c >= 'a' && c <= 'z'){</pre>
38
    return T_ID;
39
   } else if (c >= '0' && c <= '9'){</pre>
40
    char next_c = (char) inputReader.read();
41
    if (next_c >= '0' && next_c <= '9'){</pre>
42
     return T_NUM;
43
44
45
   throw new ParseException();
  }
47
```

First Stage $T \quad OP \quad t[0] \geq ' *' \wedge t[0] \leq ' +'$ $T \quad ID \quad t[0] \ge 'a' \wedge t[0] \le 'z'$ T NUM $t[0] \geq 0' \wedge t[0] \leq 9' \wedge t[0] \leq 10' \wedge t[0] < 10' \wedge t[0] \leq 10' \wedge t[0] < 10' \wedge t[0$ $t[1] \ge 0' \land t[1] \le 9'$



• Grammar-agnostic DSE

Input: "12+13" First Patl

$T[0] = T_NUM \land T[1] = T_OP \land T[2] = T_NUM \land a[4] \neq 'z'$



- DSE Second Stage
- First Path's Path Condition

• Grammar-agnostic DSE

Input: "12+13" First Patl

$\mathbf{T}[0] = \mathbf{T}_{\mathsf{NUM}} \wedge \mathbf{T}[1] = \mathbf{T}_{\mathsf{OP}} \wedge \mathbf{T}[2] = \mathbf{T}_{\mathsf{NUM}} \wedge a[4] = '\mathsf{z}'$



- DSE Second Stage
- First Path's Path Condition

Motivation Example $T[0] = T_NUM \wedge T[1] = T_OP \wedge T[2] = T_NUM \wedge a[4] = 'z'$ Application Logic PC





Motivation Example $T[0] = T_NUM \wedge T[1] = T_OP \wedge T[2] = T_NUM \wedge a[4] = 'z'$ Application Logic PC



$\begin{array}{l} \textbf{Motivation Example} \\ T[0] = T_NUM \land T[1] = T_OP \land T[2] = T_NUM \land a[4] = 'z' \\ \hline \textbf{Token PC} \end{array}$






• Grammar-agnostic DSE

Input: "12+13"

Token PC $T[0] = T_NUM \land T[1] = T_OP \land T[2] \neq T_NUM$

- Second Stage
- First Path's Path Condition

• Grammar-agnostic DSE

Input: "12+13"

Token PC $T[0] = T_NUM \land T[1] = T_OP \land T[2] \neq T_NUM$

 $T[i] \in \{T_ID, T_NUM, T_OP\}$ *i*=1

- Second Stage
- First Path's Path Condition

• Grammar-agnostic DSE

Input: "12+13" First Path's Path Condition

Token PC $T[0] = T_NUM \land T[1] = T_OP \land T[2] \neq T_NUM$

 $\bigwedge T[i] \in \{T_ID, T_NUM, T_OP\} \quad \square \land \langle T_NUM, T_OP, T_ID \rangle$ i=1

- Second Stage

 $\langle T_NUM, T_OP, T_ID \rangle$ $T_OP \quad t[0] \ge ' *' \land t[0] \le ' +'$ T [D $t[0] \ge a' \wedge t[0] \le z'$ T NUM $t[0] \geq 0' \wedge t[0] \leq 9' \wedge t[0]$ $t[1] \ge 0' \land t[1] \le 9'$

 $a[0] \ge '0' \land a[0] \le '9' \land a[1] \ge '0' \land a[1] \le '9' \land a[1]$ $a[2] \geq '*' \wedge a[2] \leq '+' \wedge$ $a[3] \ge a' \land a[3] \le z'$



 $\langle T_NUM, T_OP, T_ID \rangle$ $\mathbf{T}_{\mathbf{OP}} \quad t[0] \geq ' *' \wedge t[0] \leq ' +'$ **T ID** $t[0] \geq a' \wedge t[0] \leq z'$ T NUM $t[0] \geq 0' \wedge t[0] \leq 9' \wedge t[0] \leq 10' \wedge t[0] < 10' \wedge t[0$ $t[1] \ge 0' \land t[1] \le 9'$

 $a[0] \ge '0' \land a[0] \le '9' \land a[1] \ge '0' \land a[1] \le '9' \land a[2] \ge '*' \land a[2] \le '+' \land a[2] \ge '+' \land a[2] \le '+' \land a[2]$ $a[3] \ge a' \land a[3] \le z'$

New Input: "111+a"



• Grammar-agnostic DSE





$T[0] = T_NUM \land T[1] = T_OP \land T[2] \neq T_NUM \land T[2] = T_ID \land$ **Token PC** $a[3] \neq$ **'**7' Application Logic PC

• Grammar-agnostic DSE

Input: "11+a" Second Path's Path Condition



- Second Stage
- $T[0] = T_NUM \wedge T[1] = T_OP \wedge T[2] \neq T_NUM \wedge T[2] = T_ID \wedge$

• Grammar-agnostic DSE

Input: "11+a" Second Path's Path Condition



- Second Stage
- $T[0] = T_NUM \land T[1] = T_OP \land T[2] \neq T_NUM \land T[2] = T_ID \land$

(ID) ::= 'a' | 'b' | ... | 'y' | 'z' (Number) ::= '00' | '01' | ... | '98' | '99' ⟨Op⟩ ::= '*' | '+'

Grammar-Agnostic DSE

```
\langle Expr \rangle ::= \langle ID \rangle | \langle Number \rangle | \langle Number \rangle \langle Op \rangle \langle Expr \rangle
```

Cover all the parsing code at the 2nd iteration. Trigger the bug at the 3rd iteration

- JPF-based DSE engine
- Java Programs
- Implementation



- JPF-based DSE engine
- Java Programs
- Implementation



- Research questions
 - Effectiveness
 - Bug detection and coverage
 - Efficiency

- Benchmark
 - 19 open-source Java parsing programs
 - Il types of grammars
 - 5~128 token types

- Baseline
 - Byte-level symbolization
 - JQF: coverage-guided fuzzing [ISSTA'19]
- Setup
 - I hour analysis time

• Grammar: grammar-guided black-box fuzzing [ASE'19]

Results of Bug Finding

Name	Project	Type	GADSE	Char	JQF	Grammar
Bug 1	J2latex	NumberFormatException	143s	>1h	185s	No
Bug 2	CMMParser	NullPointerException	36s	>1h	>1h	Yes
Bug 3	CMMParser	NumberFormatException	41s	>1h	>1h	Yes
Bug 4	Jsijcc	NullPointerException	456s	>1h	>1h	No
Bug 5	Jsijcc	ClassCastException	163s	77s	>1h	No
Bug 6	Jsijcc	ClassCastException	44s	>1h	>1h	Yes

Find 6 unknown bugs, and each needs less than 8 minutes



 $N_{GADSE} - N_{CHAR}$

 $N_{\mathbf{CHAR}}$



Better than byte-level symbolization in both statement and branch coverages



Results of Coverage



Results of Coverage



GADSE achieve best results in BFS mode

Results of Coverage



94



























ACM SIGSOFT International Symposium on Software Testing and Analysis

Thank you! Q&A



