



# Relaxed Parsing of Regular Approximations of String-Embedded Languages

**Author:** Ekaterina Verbitskaia

Saint Petersburg State University  
JetBrains Programming Languages and Tools Lab

26/08/2015

# String embedding

- Embedded SQL

```
SqlCommand myCommand = new SqlCommand(  
    "SELECT * FROM table WHERE Column = @Param2",  
    myConnection);  
myCommand.Parameters.Add(myParam2);
```

- Dynamic SQL

```
IF @X = @Y  
    SET @TBL = ' #table1 '  
ELSE  
    SET @TBL = ' table2 '  
SET @S = 'SELECT x FROM' + @TBL + 'WHERE ISNULL(n,0) > 1'  
EXECUTE (@S)
```

- String-embedded code are expressions in some programming language
  - ▶ It may be necessary to support them in IDE: code highlighting, autocomplete, refactorings
  - ▶ It may be necessary to transform them: migration of legacy software on new platforms
  - ▶ It may be necessary to detect vulnerabilities in such code
  - ▶ Any other problems of programming languages can occur

# Static analysis of string-embedded code

- Performed without programm execution
- Checks that the set of properties holds for each possible expression value
- Undecidable for string-embedded code in the general case
- The set of possible expression values is over approximated and then the approximation is analysed.

# Existing tools

- PHP String Analyzer, Java String Analyzer, Alvor
  - ▶ Static analyzers for PHP, Java, and SQL embedded into Java respectively
- Kyung-Goo Doh et al.
  - ▶ Checks syntactical correctness of embedded code
- PHPStorm
  - ▶ IDE for PHP with support of HTML, CSS, JavaScript
- IntelliLang
  - ▶ PHPStorm and IDEA plugin, supports various languages
- STRANGER
  - ▶ Vulnerability detection of PHP
- Flaws
  - ▶ Limited functionality
  - ▶ Hard to extend them with new features or support new languages
  - ▶ Do not create structural representation of code

# Static analysis of string-embedded code: the scheme

- Identification of hotspots: points of interest, where the analysis is desirable
- Approximation construction
- Lexical analysis
- **Syntactic analysis**
- Semantic analysis

# Static analysis of string-embedded code: the scheme

Code: hotspot is marked

```
string res = "";  
for(i = 0; i < 1; i++)  
    res = "()" + res;  
use(res);
```

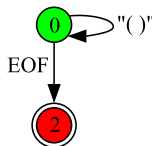
Possible values

`{"", "()", "()", ..., "()"^1}`

Regular approximation

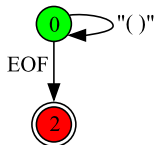
`("()")*`

Approximation

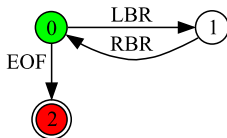


# Static analysis of string-embedded code: the scheme

## Approximation



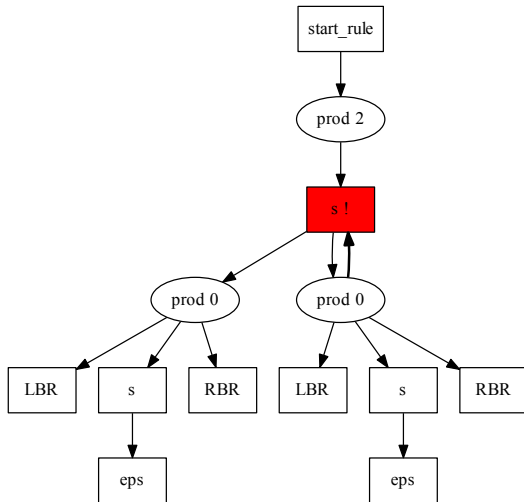
## After lexing



## Grammar

$start ::= s$   
 $s ::= LBR\ s\ RBR\ s$   
 $s ::= \epsilon$

## Parse forest





**The aim** is to develop the algorithm suitable for syntactic analysis of string-embedded code

**Tasks:**

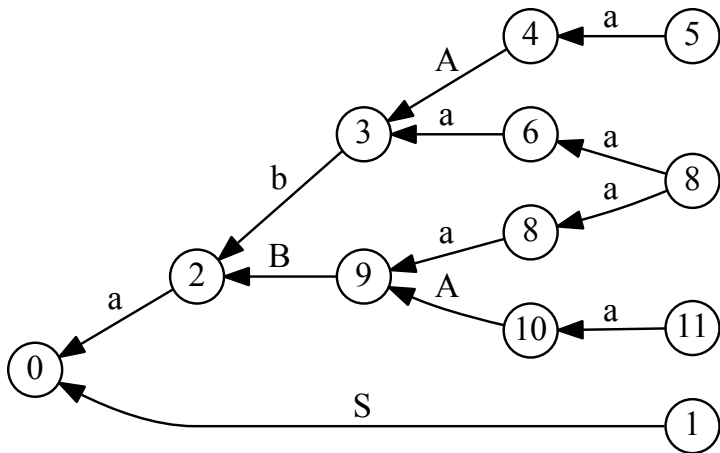
- Develop an algorithm for parsing of regular approximation of embedded code which produce a finite parse forest
- Parse forest should contain a parse tree for every correct (w.r.t. reference grammar) string accepted by the input automaton
- Incorrect strings should be omitted: no error detection
- An algorithm should not depend on the language of the host program and the language of embedded code

- **Input:** reference DCF-grammar  $G$  and DFA graph with no  $\epsilon$ -transitions over the alphabeth of terminals of  $G$
- **Output:** finite representation of the trees corresponding to all correct string accepted by input automaton

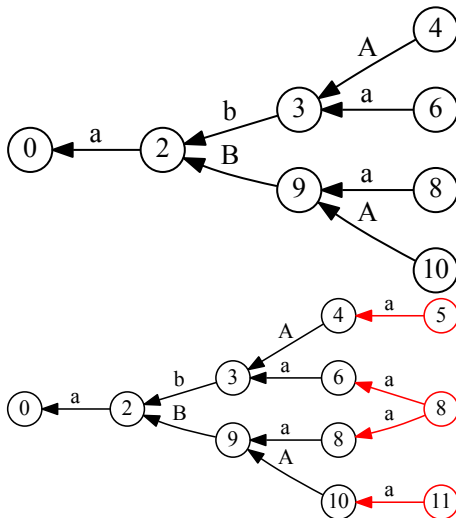
# Right-Nullled Generalized LR algorithm

- RNLGR processes context free grammars
- Uses data structures which reduce memory consumption and guarantee appropriate time of analysis
- In case when LR conflicts occur, parses in each possible way
  - ▶ Shift/Reduce conflict
  - ▶ Reduce/Reduce conflict

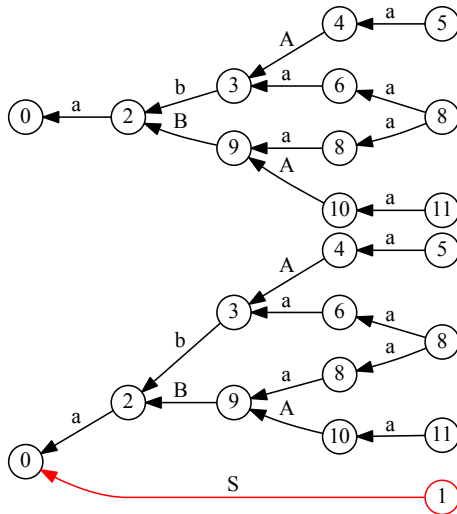
# RNGLR data structures: Graph-Structured Stack



# RNGLR operations: shift



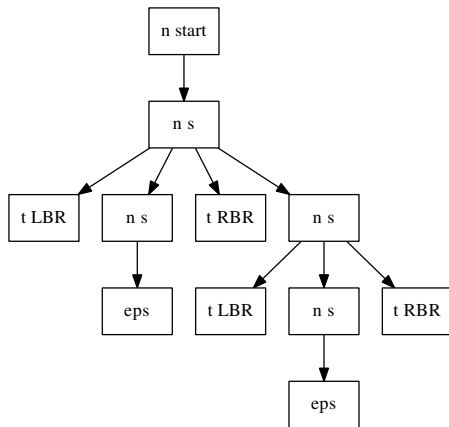
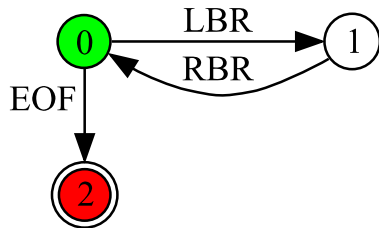
# RNGLR operations: reduce



- Traverse the automaton graph and sequentially construct GSS, similarly as in RNLRL
- The set of LR-states is associated with each vertex of input graph
- The order in which the vertices of input graph are traversed is controlled with a queue. Whenever new edge is added to GSS, its head vertex is enqueued
- The algorithm implements relaxed parsing: errors are not detected, erroneous strings are ignored

## Algorithm: correctness

*Correct tree* — derivation tree of some string accumulated along the path in the input graph





## Algorithm: correctness

### Theorem (Termination)

*Algorithm terminates for any input*

### Theorem (Correctness)

*Every tree, generated from SPPF, is correct*

### Theorem (Correctness)

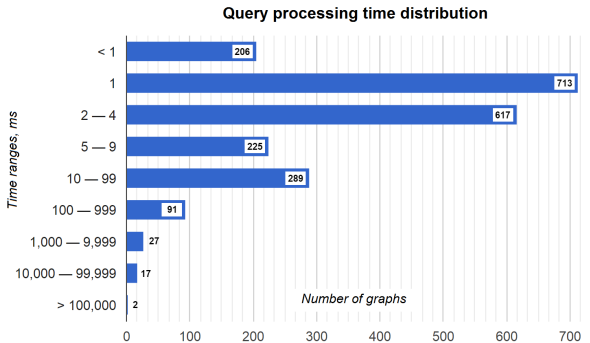
*For every path  $p$  in the inner graph, recognized w.r.t. reference grammar, a correct tree corresponding to  $p$  can be generated from SPPF*

# Implementation

- The algorithm is implemented as a part of YaccConstructor project using F# programming language
- The generator of RNLGR parse tables and data structures for GSS and SPPF are reused

# Evaluation

- The data is taken from the project of migration from MS-SQL to Oracle Server
- 2,7 lines of code, 2430 queries, 2188 successfully processed
- The number of queries which previously could not be processed because of timeout is decreased from 45 to 1



- The algorithm for parsing of regular approximation of dynamically generated string which constructs the finite representation of parse forest is developed
- Its termination and correctness are proved
- The algorithm is implemented as a part of YaccConstructor project
- The evaluation demonstrated it could be used for complex tasks