

# Scattered context grammars and generic reverse compilation

Lukáš Ďurfina

UIFS FIT VUT



## Contents

- Motivation and aims
- Ideas
- Scattered context grammar
- Application in reverse compilation

## Why?

- a lot of space for research
- usefull applications

## Probably the most usefull application

- every year milions of new unique malware programs
- families of malware with similar code and behaviour
- polymorphic and metamorphic viruses
- antivirus software, IDS

## Facts

- part of reverse engineering
- do not produce code in high level language
- find and mark patterns
- recognize specific behaviour

Give answer

YES or NO

And what was a question?

- same behaviour
- level of an accuracy

## Base

- imagine executable as word of some language  $L$
- we need mechanism for comparing these words

## Transformation

- transformational grammar
- we are [not]? able to transform one word (executable) to another



## Transformation

- transformational grammar
- we are [not]? able to transform one word (executable) to another
  - $W_A \stackrel{?}{\Rightarrow} W_B$

## Transformation

- transformational grammar
- we are [not]? able to transform one word (executable) to another
  - $W_A \stackrel{?}{\Rightarrow} W_B$
  - $W_A \Rightarrow W_{C_1}, W_B \Rightarrow W_{C_2}, W_{C_1} \stackrel{?}{=} W_{C_2}$

## Generation

- we are [not]? able to create generative grammar  $G$
- $G$  can generate “all” words (executables) with specific behaviour

## Generation

- we are [not]? able to create generative grammar  $G$
- $G$  can generate “all” words (executables) with specific behaviour
- $G$  can also parse them

## Generation

- we are [not]? able to create generative grammar  $G$
- $G$  can generate “all” words (executables) with specific behaviour
- $G$  can also parse them
  - LL-parsing

## Conditions

- has strong generative power
- can describe the problem
- is not too much complex for implementation

## Definition

$$G = (V, T, P, S)$$

$V$  is the total alphabet

$T$  is the set of terminals,  $T \subset V$

$S$  is the start symbol,  $S \in V - T$

$P$  is a finite set of productions of the form

$$(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n),$$

where  $A_1, \dots, A_n \in V - T, x_1, \dots, x_n \in V^*$

## Derivation step naturally describes problem

For  $(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n) \in P$  and

$$u = u_1 A_1 \dots u_n A_n u_{n+1}$$

$$v = u_1 x_1 \dots u_n x_n u_{n+1}$$

we write  $u \Rightarrow v[(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)]$



Strong generative power

$$\mathcal{L}(SC) = \mathcal{L}(RE)$$

Strong generative power

$$\mathcal{L}(SC) = \mathcal{L}(RE)$$

Comparison with PSCG

allow  $\epsilon$  rules

## Example

```
a = countA();  
b = countB();  
return magic(a,b);
```

```
b = countB();  
a = countA();  
return magic(a,b);
```

## Example

<code>a = countA();</code>	<code>b = countB();</code>
<code>b = countB();</code>	<code>a = countA();</code>
<code>return magic(a,b);</code>	<code>return magic(a,b);</code>

$(\langle b = \text{countB}() \rangle, \langle a = \text{countA}() \rangle) \rightarrow (\langle a = \text{countA}() \rangle, \langle b = \text{countB}() \rangle)$

## Example

<code>a = countA();</code>	<code>a = countA();</code>
<code>b = countB();</code>	<code>b = countB();</code>
<code>return magic(a,b);</code>	<code>return magic(a,b);</code>

$(\langle b = \text{countB}() \rangle, \langle a = \text{countA}() \rangle) \rightarrow (\langle a = \text{countA}() \rangle, \langle b = \text{countB}() \rangle)$

## Example

```
a = countA();  
b = countB();  
c = countC();  
return magic(a,b);
```

## Example

```
int a = countA();  
int b = countB();  
volatile int c = countC();  
return magic(a,b);
```

## Example

```
A:      int a = countA();  
B:      int b = countB();  
C:      volatile int c = countC();  
<AB>:   return magic(a,b);
```

$(A, B, C, \langle AB \rangle) \rightarrow (A, B, \epsilon, \langle AB \rangle)$



## Example

```
A:      int a = countA();  
B:      int b = countB();  
C:      volatile int c = countC();  
⟨AB⟩:   return magic(a,b);
```

$$(A, B, C, \langle AB \rangle) \rightarrow (\textcolor{red}{A}, \textcolor{green}{B}, \epsilon, \langle \textcolor{red}{A}\textcolor{green}{B} \rangle)$$

## Example

```
int a = countA();  
int b = countB();  
return magic(a,b);
```

## Definition

$$G = (V, T, P, I)$$

$V$  is the total vocabulary

$T$  is the set of terminals (or the *output vocabulary*),  $T \subset V$

$P$  is a finite set of productions of scattered context productions

$I$  is the *input vocabulary*,  $I \subset V$

## Derivation step

$T(G, K)$ : transformation  $T$  that  $G$  defines from  $K \subseteq I^*$

$$T(G, K) = \{(x, y) : x \Rightarrow_G^* y, x \in K, y \in T^*\}$$

## Example

$$G = (V, T, P, I)$$

$$V = \{A, B, C, a, b, c\}$$

$$T = \{a, b, c\}$$

$$I = \{A, B, C\}$$

$$P = \{(A, B, C) \rightarrow (a, bb, c)\}$$

We can take input sentence  $AABBCC$ :

$$AABBCC \Rightarrow_G aABbbcC \Rightarrow_G aabbbbcc$$

## Action blocks

- granularity
- connections

## Example

INC ax	INC bx
INC bx	INC ax
MUL bx	MUL bx
ADD bx, ax	MOV ebx, eax
PUSH eax	PUSH eax

## Example

$\langle INCA \rangle$

$\langle INCB \rangle$

$\langle MULB \rangle$

$\langle ADDBA \rangle$

$\langle PUSHA \rangle$

$\langle INCB \rangle$

$\langle INCA \rangle$

$\langle MULB \rangle$

$\langle MOVBA \rangle$

$\langle PUSHA \rangle$

## Example

$\langle INCA \rangle$	$\langle INCB \rangle$
$\langle INCB \rangle$	$\langle INCA \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle MOVBA \rangle$
$\langle PUSHA \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$



## Example

$\langle INCA \rangle$	$\langle INCB \rangle$
$\langle INCB \rangle$	$\langle INCA \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle MOVBA \rangle$
$\langle PUSHA \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle MOVBA \rangle$
$\langle PUSHA \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle MOVBA \rangle$
$\langle PUSHA \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle \textcolor{red}{MULB} \rangle$
$\langle ADDBA \rangle$	$\langle \textcolor{brown}{MOVBA} \rangle$
$\langle PUSHA \rangle$	$\langle \textcolor{blue}{PUSHA} \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle \textcolor{red}{MULB} \rangle, \langle \textcolor{brown}{MOVBA} \rangle, \langle \textcolor{blue}{PUSHA} \rangle) \rightarrow (\langle \textcolor{red}{MULB} \rangle, \epsilon, \langle \textcolor{blue}{PUSHA} \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle PUSHBA \rangle$
$\langle PUSHBA \rangle$	

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHBA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHBA \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle ADDBA \rangle$	$\langle PUSHA \rangle$
$\langle PUSHA \rangle$	

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$
$$(\langle MULB \rangle, \langle ADDBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle \textcolor{red}{MULB} \rangle$	$\langle MULB \rangle$
$\langle \textcolor{brown}{ADDBA} \rangle$	$\langle PUSH A \rangle$
$\langle \textcolor{blue}{PUSHA} \rangle$	

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSH A \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSH A \rangle)$$
$$(\langle \textcolor{red}{MULB} \rangle, \langle \textcolor{brown}{ADDBA} \rangle, \langle \textcolor{blue}{PUSHA} \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSH A \rangle)$$

## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle \textcolor{red}{MULB} \rangle$	$\langle MULB \rangle$
$\langle \textcolor{blue}{PUSHA} \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$
$$(\langle MULB \rangle, \langle ADDBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle \textcolor{red}{MULB} \rangle, \epsilon, \langle \textcolor{blue}{PUSHA} \rangle)$$



## Example

$\langle INCA \rangle$	$\langle INCA \rangle$
$\langle INCB \rangle$	$\langle INCB \rangle$
$\langle MULB \rangle$	$\langle MULB \rangle$
$\langle PUSHA \rangle$	$\langle PUSHA \rangle$

Rules:

$$(\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle)$$
$$(\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$
$$(\langle MULB \rangle, \langle ADDBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle)$$

## Example

$$G = (V, T, P, I)$$

$$V = \{ \langle INCA \rangle, \langle INCB \rangle, \langle MULB \rangle, \langle ADDBA \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle, \\ INCA, INCB, MULB, ADDBA, MOVBA, PUSHA \}$$

$$T = \{ INCA, INCB, MULB, ADDBA, MOVBA, PUSHA \}$$

$$I = \{ \langle INCA \rangle, \langle INCB \rangle, \langle MULB \rangle, \langle ADDBA \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle \}$$

$$P = \{ (\langle INCB \rangle, \langle INCA \rangle) \rightarrow (\langle INCA \rangle, \langle INCB \rangle), \\ (\langle MULB \rangle, \langle MOVBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle), \\ (\langle MULB \rangle, \langle ADDBA \rangle, \langle PUSHA \rangle) \rightarrow (\langle MULB \rangle, \epsilon, \langle PUSHA \rangle), \\ (\langle INCA \rangle) \rightarrow (INCA), \dots \}$$

Thank you for your attention



A. Meduna and J. Techet.

*Scattered Context Grammars and their Applications.*

WIT Press, 2010.



P. Szor.

*Počítačové viry - analýza útoku a obrana.*

Zoner Press, 2006.