

39th International Conference on Current Trends  
in Theory and Practice of Computer Science  
Špindleruv Mlýn, Czech Republic

## Recursive descent parsing for grammars with contexts

Mikhail Barash

Ph.D. student,  
Department of Mathematics and Statistics, University of Turku  
Turku Centre for Computer Science (TUCS)  
Finland

January 26–31, 2013

# What a recursive descent is?

- subclass of grammars allowing recursive descent:  
LL( $k$ )-grammars
  - $k$  look-ahead symbols

# What a recursive descent is?

- subclass of grammars allowing recursive descent:  
LL( $k$ )-grammars
  - $k$  look-ahead symbols
- LL(1) context-free grammar generating  $\{ a^n b^n \mid n \geq 0 \}$ :

$$S \rightarrow aSb \mid \varepsilon$$

<pre>S() {     if (look-ahead is "a") {         a(); S(); b();     } }</pre>	<pre>a() {     if (current symbol is "a")         advance position by 1     else error }</pre>
--	--

# What a recursive descent is?

- subclass of grammars allowing recursive descent:  
LL( $k$ )-grammars
  - $k$  look-ahead symbols
- LL(1) context-free grammar generating  $\{ a^n b^n \mid n \geq 0 \}$ :

$$S \rightarrow aSb \mid \varepsilon$$

<pre>S() {     if (look-ahead is "a") {         a(); S(); b();     } }</pre>	<pre>a() {     if (current symbol is "a")         advance position by 1     else error }</pre>
--	--

- First compilers for Pascal are recursive descent parsers.
  - Implemented "by hand" (N. Wirth, 1970).
  - Program code can be easily generated automatically.

## Conjunctive grammars (Okhotin, 2001)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k, \quad \alpha_j \in (\Sigma \cup N)^*$$

## Conjunctive grammars (Okhotin, 2001)

$A \rightarrow \alpha_1 \& \dots \& \alpha_k, \quad \alpha_j \in (\Sigma \cup N)^*$

- “a string  $w$  has property  $A \iff w$  has all the properties  $\alpha_1, \dots, \alpha_k$ ”

# Extension of CFGs with Boolean operations

## Conjunctive grammars (Okhotin, 2001)

$A \rightarrow \alpha_1 \& \dots \& \alpha_k, \quad \alpha_j \in (\Sigma \cup N)^*$

- “a string  $w$  has property  $A \iff w$  has all the properties  $\alpha_1, \dots, \alpha_k$ ”

$\{ a^n b^n c^n \mid n \geq 0 \}$

$S \rightarrow AB \& DC$

$A \rightarrow aA \mid \varepsilon$

$B \rightarrow bBc \mid \varepsilon$

$C \rightarrow cC \mid \varepsilon$

$D \rightarrow aDb \mid \varepsilon$

## Conjunctive grammars (Okhotin, 2001)

$A \rightarrow \alpha_1 \& \dots \& \alpha_k, \quad \alpha_j \in (\Sigma \cup N)^*$

- “a string  $w$  has property  $A \iff w$  has all the properties  $\alpha_1, \dots, \alpha_k$ ”

$\{ a^n b^n c^n \mid n \geq 0 \}$

$S \rightarrow AB \& DC$

$A \rightarrow aA \mid \varepsilon$

$B \rightarrow bBc \mid \varepsilon$

$C \rightarrow cC \mid \varepsilon$

$D \rightarrow aDb \mid \varepsilon$

- non-context-free languages generated
- complexity of basic parsing algorithms preserved
- nontrivial properties of subclasses



$$A \rightarrow \alpha_1 \& \dots \& \alpha_k$$

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m$$

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = u w v$ :

- each  $\alpha_i$  defines  $w$

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = u w v$ :

- each  $\alpha_i$  defines  $w$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = u w v$ :

- each  $\alpha_i$  defines  $w$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )
- each  $\gamma_i$  defines  $wv$  (the extended right context of  $w$ )

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = u w v$ :

- each  $\alpha_i$  defines  $w$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )
- each  $\gamma_i$  defines  $wv$  (the extended right context of  $w$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle w \rangle v]$  : “a string  $w$  written in right context  $v$  has the property  $X$ ”

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = u w v$ :

- each  $\alpha_i$  defines  $w$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )
- each  $\gamma_i$  defines  $wv$  (the extended right context of  $w$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle w \rangle v]$  : “a string  $w$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2^n} d\}$$



# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = uv$ :

- each  $\alpha_i$  defines  $u$
- each  $\beta_i$  defines  $v$  (the right context of  $u$ )
- each  $\gamma_i$  defines  $uv$  (the extended right context of  $u$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle u \rangle v]$  : “a string  $u$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2^n} d\}$$

$$S \rightarrow Ac \mid Bd$$

$$A \rightarrow aAb \mid \varepsilon$$

$$B \rightarrow aBbb \mid \varepsilon$$

$$L(A) = \{a^n b^n\}$$

$$L(B) = \{a^n b^{2^n}\}$$

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = uv$ :

- each  $\alpha_i$  defines  $u$
- each  $\beta_i$  defines  $v$  (the right context of  $u$ )
- each  $\gamma_i$  defines  $uv$  (the extended right context of  $u$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle u \rangle v]$  : “a string  $u$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2n} d\}$$

$$S \rightarrow Ac \mid Bd$$

$$A \rightarrow aAb \mid \varepsilon$$

$$B \rightarrow aBbb \mid \varepsilon$$

$$L(A) = \{a^n b^n\}$$

$$L(B) = \{a^n b^{2n}\}$$

- the language is **not standard LL( $k$ )** for any  $k$

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = uv$ :

- each  $\alpha_i$  defines  $u$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )
- each  $\gamma_i$  defines  $wv$  (the extended right context of  $w$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle w \rangle v]$  : “a string  $w$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2n} d\}$$

$$S \rightarrow Ac \mid Bd$$

$$A \rightarrow aAb \mid \varepsilon$$

$$B \rightarrow aBbb \mid \varepsilon$$

$$L(A) = \{a^n b^n\}$$

$$L(B) = \{a^n b^{2n}\}$$

- the language is **not standard LL( $k$ )** for any  $k$
- way out: use **right contexts to peek the last symbol** of the string

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = uv$ :

- each  $\alpha_i$  defines  $u$
- each  $\beta_i$  defines  $v$  (the right context of  $u$ )
- each  $\gamma_i$  defines  $uv$  (the extended right context of  $u$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle u \rangle v]$  : “a string  $u$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2n} d\}$$

$$S \rightarrow Ac \mid Bd$$

$$A \rightarrow aAb \mid \varepsilon$$

$$B \rightarrow aBbb \mid \varepsilon$$

$$X \rightarrow aX \mid bX \mid \varepsilon$$

$$L(A) = \{a^n b^n\}$$

$$L(B) = \{a^n b^{2n}\}$$

$$L(X) = (a \mid b)^*$$

- the language is **not standard LL( $k$ )** for any  $k$
- way out: use **right contexts to peek the last symbol** of the string

# Grammars with contexts (B, Okhotin, 2012)

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleright \beta_1 \& \dots \& \triangleright \beta_m \& \triangleright \gamma_1 \& \dots \& \triangleright \gamma_n$$

Given a string  $x = uv$ :

- each  $\alpha_i$  defines  $u$
- each  $\beta_i$  defines  $v$  (the right context of  $w$ )
- each  $\gamma_i$  defines  $wv$  (the extended right context of  $w$ )

Semantics by *logical deduction* of elementary propositions

$[X, \langle w \rangle v]$  : “a string  $w$  written in right context  $v$  has the property  $X$ ”

$$\{a^n b^n c\} \cup \{a^n b^{2n} d\}$$

$$S \rightarrow \triangleright Xc \& Ac \mid \triangleright Xd \& Bd$$

$$A \rightarrow aAb \mid \varepsilon$$

$$L(A) = \{a^n b^n\}$$

$$B \rightarrow aBbb \mid \varepsilon$$

$$L(B) = \{a^n b^{2n}\}$$

$$X \rightarrow aX \mid bX \mid \varepsilon$$

$$L(X) = (a \mid b)^*$$

- the language is **not standard LL( $k$ )** for any  $k$
- way out: use **right contexts to peek the last symbol** of the string

Recursive descent parser:

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules



Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Theory:

- **mathematically sound definition** of right contexts (*"look-ahead"*) within recursive descent

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Theory:

- **mathematically sound definition** of right contexts (*"look-ahead"*) within recursive descent
- **higher expressive power**, but still linear time

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Theory:

- **mathematically sound definition** of right contexts (*"look-ahead"*) within recursive descent
- **higher expressive power**, but still linear time
- parsing algorithm **proven correct**

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Theory:

- **mathematically sound definition** of right contexts (*"look-ahead"*) within recursive descent
- **higher expressive power**, but still linear time
- parsing algorithm **proven correct**

Practice:

- implemented as a **prototype software**

Recursive descent parser:

- **conjunction in rules**: scan the substring multiple times
- **right contexts** ( $\triangleright, \trianglerighteq$ ): to choose suitable alternatives of rules
- **backtracking**: to go over the alternatives
- **memoization**: linear time instead of exponential

Theory:

- **mathematically sound definition** of right contexts (*"look-ahead"*) within recursive descent
- **higher expressive power**, but still linear time
- parsing algorithm **proven correct**

Practice:

- implemented as a **prototype software**

