

Top-Down Syntax Analysis

Kristoffer H. Rose Eva Rose
krisrose@cs.nyu.edu

Compiler Construction
CSCI-GA.2130-001/Spring 2014
NYU Courant Institute

February 10, 2014

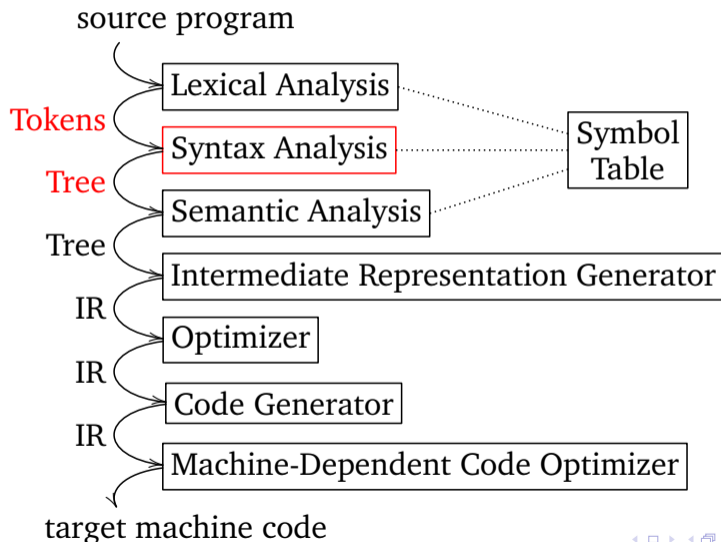


Outline

- 1 Parsers
- 2 Context-Free Grammars
- 3 Grammar Transforms & HACS
- 4 Top-Down Parsing
- 5 HACS Parsing



Second compilation phase



Example

```
position = initial + rate * 60
```

scanned into list of *tokens*:

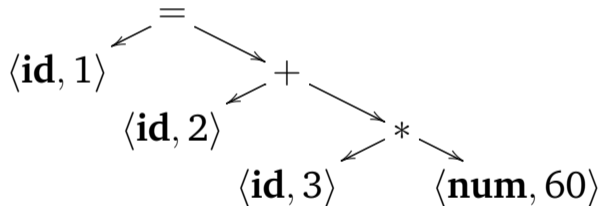
$\langle \mathbf{id}, 1 \rangle$ $\langle = \rangle$ $\langle \mathbf{id}, 2 \rangle$ $\langle + \rangle$ $\langle \mathbf{id}, 3 \rangle$ $\langle * \rangle$ $\langle \mathbf{num}, 60 \rangle$

1	position
2	initial
3	rate



Example

parsed into parse/abstract syntax tree:



according to *precedence rules...*



Some syntax rules.

Some tree-structuring syntax rules:

- ▶ block indentation rules (Python)
- ▶ block delimiters like `{...}` (Java, C)
- ▶ grouping rules like `(...)` (most languages)
- ▶ built-in algebraic precedence rules (most languages)
- ▶ statements vs expressions ...



Outline

- 1 **Parsers**
- 2 Context-Free Grammars
- 3 Grammar Transforms & HACS
- 4 Top-Down Parsing
- 5 HACS Parsing



A Left Recursive Grammar

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned} \tag{4.1}$$

Bottom-Up Construction



A Left Recursive Grammar

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned} \tag{4.1}$$

Bottom-Up Construction



A Non-Left-Recursive Grammar

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned} \tag{4.2}$$

Top-Down Construction



A Non-Left-Recursive Grammar

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned} \tag{4.2}$$

Top-Down Construction



An Abstract Grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid \mathbf{id} \quad (4.3)$$

Good for *already constructed* trees!



An Abstract Grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid \mathbf{id} \quad (4.3)$$

Good for *already constructed* trees!



Errors

- ▶ Lexical.
- ▶ Syntactic.
- ▶ Semantic.
- ▶ Logical.

Detect early? Recover? Avoid repetition?

- ▶ Panic! (for a while, “skip to ;”)
- ▶ Phrase-level.
- ▶ Error Productions.
- ▶ Global Correction.



Errors

- ▶ Lexical.
- ▶ Syntactic.
- ▶ Semantic.
- ▶ Logical.

Detect early? Recover? Avoid repetition?

- ▶ Panic! (for a while, “skip to ;”)
- ▶ Phrase-level.
- ▶ Error Productions.
- ▶ Global Correction.



Errors

- ▶ Lexical.
- ▶ Syntactic.
- ▶ Semantic.
- ▶ Logical.

Detect early? Recover? Avoid repetition?

- ▶ Panic! (for a while, “skip to ;”)
- ▶ Phrase-level.
- ▶ Error Productions.
- ▶ Global Correction.



Errors

- ▶ Lexical.
- ▶ Syntactic.
- ▶ Semantic.
- ▶ Logical.

Detect early? Recover? Avoid repetition?

- ▶ Panic! (for a while, “skip to ;”)
- ▶ Phrase-level.
- ▶ Error Productions.
- ▶ Global Correction.



Errors

- ▶ Lexical.
- ▶ Syntactic.
- ▶ Semantic.
- ▶ Logical.

Detect early? Recover? Avoid repetition?

- ▶ Panic! (for a while, “skip to ;”)
- ▶ Phrase-level.
- ▶ Error Productions.
- ▶ Global Correction.



Outline

- 1 Parsers
- 2 Context-Free Grammars**
- 3 Grammar Transforms & HACS
- 4 Top-Down Parsing
- 5 HACS Parsing



Definition

Context-free Grammar

$$G = (V, \Sigma, R, S)$$

where

- ▶ ' $V \in V$ ' – finite set of *non-terminals*
- ▶ ' Σ ' – finite set of *terminals* (or *alphabet*, where $\Sigma \cap V = \emptyset$)
- ▶ ' R ' – finite relation of *productions, or rewrite rules*
($R \subseteq V \times (V \cup \Sigma)^*$)
- ▶ ' S ' – *start symbol* ($S \in V$)



Fig. 4.2

$$\text{expression} \rightarrow \text{expression} + \text{term}$$
$$\text{expression} \rightarrow \text{expression} - \text{term}$$
$$\text{expression} \rightarrow \text{term}$$
$$\text{term} \rightarrow \text{term} * \text{factor}$$
$$\text{term} \rightarrow \text{term} / \text{factor}$$
$$\text{term} \rightarrow \text{factor}$$
$$\text{factor} \rightarrow (\text{expression})$$
$$\text{factor} \rightarrow \mathbf{id}$$


Derivations

“Unfolding grammar until only tokens left.”

Ambiguity = “More than one (left-/rightmost) derivation.”



Derivations

“Unfolding grammar until only tokens left.”

Ambiguity = “More than one (left-/rightmost) derivation.”



Parse Tree

“Tracing unfolding grammar.”

- ▶ Each unfolded nonterminal leaves a node.
- ▶ Each token a leaf.



Parse Tree

“Tracing unfolding grammar.”

- ▶ Each unfolded nonterminal leaves a node.
- ▶ Each token a leaf.



CFG vs RE

“CFG can simulate any RE.”



CFG vs RE

“CFG can simulate any RE.”

- ▶ RE state = nonterminal.
- ▶ RE transition = production $S \rightarrow tS'$.
- ▶ RE accept = empty production.
- ▶ RE start state = start symbol.



CFG vs RE

“CFG can simulate any RE.”

- ▶ RE state = nonterminal.
- ▶ RE transition = production $S \rightarrow tS'$.
- ▶ RE accept = empty production.
- ▶ RE start state = start symbol.



CFG vs RE

“CFG can simulate any RE.”

- ▶ RE state = nonterminal.
- ▶ RE transition = production $S \rightarrow tS'$.
- ▶ RE accept = empty production.
- ▶ RE start state = start symbol.



CFG vs RE

“CFG can simulate any RE.”

- ▶ RE state = nonterminal.
- ▶ RE transition = production $S \rightarrow tS'$.
- ▶ RE accept = empty production.
- ▶ RE start state = start symbol.



Outline

- 1 Parsers
- 2 Context-Free Grammars
- 3 Grammar Transforms & HACS**
- 4 Top-Down Parsing
- 5 HACS Parsing



HACS Encoding

```
sort  $E$  |  $[[\langle E \rangle + \langle T \rangle]]$  |  $[[\langle T \rangle]]$  ;  
sort  $T$  |  $[[\langle T \rangle * \langle F \rangle]]$  |  $[[\langle F \rangle]]$  ;  
sort  $F$  |  $[[ (\langle E \rangle) ]]$  |  $[[\langle Id \rangle]]$  ;  
token  $Id$  |  $[a-z]^+$  ;
```

- ▶ Explicit " $[[]]$ " for *concrete syntax*.
- ▶ Explicit " $\langle \rangle$ " for *nonterminal and terminal references*.
- ▶ "sort" means "kind of parse tree node."
- ▶ Immediate left recursion is alright.
- ▶ (All top level alternatives introduced with "|".)



HACS Encoding

```

sort  $E$  |  $[[\langle E \rangle + \langle T \rangle]]$  |  $[[\langle T \rangle]]$  ;
sort  $T$  |  $[[\langle T \rangle * \langle F \rangle]]$  |  $[[\langle F \rangle]]$  ;
sort  $F$  |  $[[ (\langle E \rangle) ]]$  |  $[[\langle Id \rangle]]$  ;
token  $Id$  |  $[a-z]^+$  ;

```

- ▶ Explicit " $[[[]]]$ " for *concrete syntax*.
- ▶ Explicit " $\langle \rangle$ " for *nonterminal and terminal references*.
- ▶ "sort" means "kind of parse tree node."
- ▶ Immediate left recursion is alright.
- ▶ (All top level alternatives introduced with "|".)



Eliminating Immediate Left Recursion

$$A \rightarrow A\alpha_1 \mid \cdots \mid A\alpha_m \mid \beta_1 \mid \cdots \mid \beta_n \quad \Longrightarrow \quad \begin{array}{l} A \rightarrow \beta_1 A' \mid \cdots \mid \beta_n A' \\ A' \rightarrow \alpha_1 A' \mid \cdots \mid \alpha_m A' \mid \epsilon \end{array}$$

```
sort E | [[⟨T⟩ ⟨E2⟩]] ; sort E2 | [[+ ⟨T⟩ ⟨E2⟩]] | [[]];
sort T | [[⟨F⟩ ⟨T2⟩]] ; sort T2 | [[* ⟨F⟩ ⟨T2⟩]] | [[]];
sort F | [[(⟨E#1⟩)]] | [[⟨Id⟩]] ;
token Id | [a-z]+ ;
```



Eliminating Immediate Left Recursion

$$A \rightarrow A\alpha_1 \mid \cdots \mid A\alpha_m \mid \beta_1 \mid \cdots \mid \beta_n \quad \Longrightarrow \quad \begin{array}{l} A \rightarrow \beta_1 A' \mid \cdots \mid \beta_n A' \\ A' \rightarrow \alpha_1 A' \mid \cdots \mid \alpha_m A' \mid \epsilon \end{array}$$

sort $E \mid \llbracket \langle T \rangle \langle E2 \rangle \rrbracket$; **sort** $E2 \mid \llbracket + \langle T \rangle \langle E2 \rangle \rrbracket \mid \llbracket \rrbracket$;
sort $T \mid \llbracket \langle F \rangle \langle T2 \rangle \rrbracket$; **sort** $T2 \mid \llbracket * \langle F \rangle \langle T2 \rangle \rrbracket \mid \llbracket \rrbracket$;
sort $F \mid \llbracket (\langle E\#1 \rangle) \rrbracket \mid \llbracket \langle Id \rangle \rrbracket$;
token $Id \mid [a-z]^+$;



Eliminating General Left Recursion

But what about

$$S \rightarrow A a$$

$$A \rightarrow A c \mid S d \mid \epsilon$$

Inline, then eliminate as before:

$$\begin{array}{l}
 S \rightarrow A a \\
 A \rightarrow A c \mid A a d \mid \epsilon
 \end{array}
 \implies
 \begin{array}{l}
 S \rightarrow A a \\
 A \rightarrow A' \\
 A' \rightarrow c A' \mid a d A' \mid \epsilon
 \end{array}$$



Eliminating General Left Recursion

But what about

$$S \rightarrow A a$$

$$A \rightarrow A c \mid S d \mid \epsilon$$

Inline, then eliminate as before:

$$\begin{array}{l}
 S \rightarrow A a \\
 A \rightarrow A c \mid A a d \mid \epsilon
 \end{array}
 \Longrightarrow
 \begin{array}{l}
 S \rightarrow A a \\
 A \rightarrow A' \\
 A' \rightarrow c A' \mid a d A' \mid \epsilon
 \end{array}$$



Hint—Dumbing It Down

Sometimes you have

$$S \rightarrow AB \mid aC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$

What to do with an a ?

$$S \rightarrow AB \mid AC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$



Hint—Dumbing It Down

Sometimes you have

$$S \rightarrow AB \mid aC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$

What to do with an a ?

$$S \rightarrow AB \mid AC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$



Hint—Dumbing It Down

Sometimes you have

$$S \rightarrow AB \mid aC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$

What to do with an a ?

$$S \rightarrow AB \mid \color{red}A C$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$



Hint—Dumbing It Down

Sometimes you have

$$S \rightarrow AB \mid aC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$

What to do with an a ?

$$S \rightarrow AB \mid AC$$

$$A \rightarrow a \mid Ab$$

$$C \rightarrow c$$



Hint—HACS Encoding with Precedence

```

sort  $E$  | [[ $\langle E \rangle + \langle E@1 \rangle$ ]]
        | [[ $\langle E@1 \rangle * \langle E@2 \rangle$ ]]@1
        | sugar [[( $\langle E\#1 \rangle$ )]@2  $\rightarrow E\#1$ 
        | [[ $\langle Id \rangle$ ]] ;
  
```

- ▶ Single "sort" captures *abstract syntax tree*.
- ▶ *Precedence* marker "@" n disambiguates parse.
- ▶ "sugar" specifies concrete-only syntax.



Hint—HACS Encoding with Precedence

```

sort  $E$  | [[ $\langle E \rangle + \langle E@1 \rangle$ ]]
          | [[ $\langle E@1 \rangle * \langle E@2 \rangle$ ]]@1
          | sugar [[( $\langle E\#1 \rangle$ )]@2  $\rightarrow E\#1$ 
          | [[ $\langle Id \rangle$ ]] ;
  
```

- ▶ Single "sort" captures *abstract syntax tree*.
- ▶ *Precedence* marker "@*n*" disambiguates parse.
- ▶ "sugar" specifies concrete-only syntax.



Outline

- 1 Parsers
- 2 Context-Free Grammars
- 3 Grammar Transforms & HACS
- 4 Top-Down Parsing**
- 5 HACS Parsing



A Non-Left-Recursive Grammar

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned} \tag{4.2}$$



FIRST

- 1 If X is *terminal*, $\text{FIRST}(X) = \{X\}$.
- 2 If X is *non-terminal* with $X \rightarrow Y_1 \dots Y_k$ ($k > 0$), then $\text{FIRST}(X)$
 - ▶ includes every $a \in \text{FIRST}(Y_i)$ where $Y_1 \dots Y_{i-1} \Rightarrow \epsilon$, and
 - ▶ includes ϵ if $\forall i = 1 \dots k : \epsilon \in \text{FIRST}(Y_i)$.
- 3 if X is *non-terminal* with $X \rightarrow \epsilon$ then $\epsilon \in \text{FIRST}(X)$.



FOLLOW

- 1 $\$ \in \text{FOLLOW}(S)$.
- 2 For every production context $A \rightarrow \alpha B \beta$,
 $(\text{FIRST}(\beta) \setminus \{\epsilon\}) \subseteq \text{FOLLOW}(B)$; and
- 3 For every production context $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ where $\epsilon \in \text{FIRST}(B)$, $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$.



LL(1)—Left (to right) Leftmost (derivation) of 1 (input symbol)

Three requirements for productions $A \rightarrow \alpha \mid \beta$:

- 1 $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint (including for ϵ).
- 2 If $\epsilon \in \text{FIRST}(\alpha)$ then $\text{FIRST}(\beta)$ and $\text{FOLLOW}(A)$ are disjoint.
- 3 If $\epsilon \in \text{FIRST}(\beta)$ then $\text{FIRST}(\alpha)$ and $\text{FOLLOW}(A)$ are disjoint.



Predictive Parsing Table

Input: Grammar G .

Output: Parsing table M .

Method: For each production $A \rightarrow \alpha$ in the grammar:

- ① For each $a \in \text{FIRST}(\alpha)$ add $A \rightarrow \alpha$ to $M[A, a]$.
- ② If $\epsilon \in \text{FIRST}(\alpha)$ then for each $b \in \text{FOLLOW}(A)$ add $A \rightarrow \alpha$ to $M[A, b]$. (If $\epsilon \in \text{FIRST}(\alpha)$ and $\$ \in \text{FOLLOW}(A)$ then add $A \rightarrow \alpha$ to $M[A, \$]$.)

Any entries $M[A, a]$ that have no content are set to **error**.

$$E \rightarrow TE' \quad E' \rightarrow +TE' \mid \epsilon \quad T \rightarrow FT' \quad T' \rightarrow *FT' \mid \epsilon \quad F \rightarrow$$



Predictive Parsing Table

Input: Grammar G .

Output: Parsing table M .

Method: For each production $A \rightarrow \alpha$ in the grammar:

- ① For each $a \in \text{FIRST}(\alpha)$ add $A \rightarrow \alpha$ to $M[A, a]$.
- ② If $\epsilon \in \text{FIRST}(\alpha)$ then for each $b \in \text{FOLLOW}(A)$ add $A \rightarrow \alpha$ to $M[A, b]$. (If $\epsilon \in \text{FIRST}(\alpha)$ and $\$ \in \text{FOLLOW}(A)$ then add $A \rightarrow \alpha$ to $M[A, \$]$.)

Any entries $M[A, a]$ that have no content are set to **error**.

$$E \rightarrow T E' \quad E' \rightarrow + T E' \mid \epsilon \quad T \rightarrow F T' \quad T' \rightarrow * F T' \mid \epsilon \quad F \rightarrow$$



Parser Table

Non-terminal	<i>InputSymbol</i>					
	id	+	*	()	\$
<i>E</i>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<i>E'</i>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<i>T</i>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<i>T'</i>		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<i>F</i>	$F \rightarrow \mathbf{id}$			$F \rightarrow (E)$		



Outline

- 1 Parsers
- 2 Context-Free Grammars
- 3 Grammar Transforms & HACS
- 4 Top-Down Parsing
- 5 HACS Parsing**



HACS Expression Parser

```
1  /* 2. SYNTAX ANALYSIS. */  
  
2  sort Stat  |  [[ <Name> := <Exp> ; ] ] |  [[ { <Stat*> } ] ] ;  
  
3  sort Exp   |  [[ <Exp@1> + <Exp@2> ] ]@1  
4             |  [[ <Exp@2> * <Exp@3> ] ]@2  
5             |  [[ <Int> ] ]@3  
6             |  [[ <Float> ] ]@3  
7             |  [[ <Name> ] ]@3  
8             |  sugar [[ (<Exp#>) ] ]@3 → Exp# ;  
  
9  sort Name  |  symbol [[ <Id> ] ] ;
```



HACS Parser Generation Tricks

- ▶ Automatically does **immediate left recursion elimination**.
- ▶ Unusual brackets `[]⟨⟩...`
- ▶ “sort” covers multiple non-terminals.
- ▶ Precedence handled automatically (the @ markers).
- ▶ Sugar.
- ▶ (Names – that’s next week.)

Oh, and it is just a front-end for JavaCC.



HACS Parser Generation Tricks

- ▶ Automatically does **immediate left recursion elimination**.
- ▶ Unusual brackets `[[]]` `< >`...
- ▶ “sort” covers multiple non-terminals.
- ▶ Precedence handled automatically (the `@` markers).
- ▶ Sugar.
- ▶ (Names – that’s next week.)

Oh, and it is just a front-end for JavaCC.



Questions?

krisrose@cs.nyu.edu

