

# Syntax Analysis

**Martin Sulzmann**

# Syntax Analysis

## Objective

Recognize individual tokens as sentences of a language (beyond regular languages).

## Example 1 (OK)

Program text `x := x + 2` yields token stream

*(ident "x") assign (ident "x") plus (const 2)*

## Example 2 (Not OK)

Program text `x x := + 2` yields token stream

*(ident "x") (ident "x") assign plus (const 2)*

We will discuss later if the 'OK' example shall be semantically valid.

# Approach: Parsing with CFG

## Context-Free Grammar (CFG)

Language described in terms of a context-free grammar:

$$\textit{Command} \rightarrow \textit{ident assign Expression} \mid \dots$$
$$\textit{Expression} \rightarrow \textit{const} \mid \textit{ident} \mid \textit{Expression plus Expression}$$

*ident*, *const*, *plus* being tokens where attributes are omitted.

Note: CFG uses uppercase for non-terminal and lowercase for terminal symbols. Our parsing tool OCamlyacc just uses it the other way around!

## Parsing yields Abstract Syntax Tree (AST)

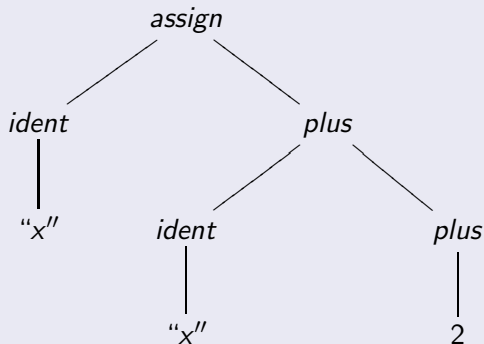
Check if stream of tokens is *valid* and produce AST.

# Parsing with CFG Example

## Token Stream

*(ident "x") assign (ident "x") plus (const 2)*

## AST



# Types of Parsing

## Top-Down

- Recursive descent
- From root to leaves
- “ANTLR” style

## Bottom-Up

- From leaves to root
- “yacc” style

## Others

- Parser combinators
- Derivative-based parser
- ...

# Error Handling

## Types of Errors

lexical	misspelled keyword	<code>while cond ...</code>
syntactic	unbalanced parentheses	<code>x:= x+1)</code>
semantic	incompatible types	<code>x:= x + True</code>
logical	infinite loop	<code>while True ...</code>

Provide useful feedback to user. Mostly ignored here.

# Context Free Grammar (CFG)

## Definition

A CFG  $G$  is a 4-tuple  $(V_t, V_n, S, P)$  where

- 1  $V_t$  is the set of *terminal* symbols (commonly the set of tokens).
- 2  $V_n$  is the set of *non-terminal symbols*.
- 3  $S$  is a distinguished non-terminal symbol call the *start* symbol.
- 4  $P$  is a finite set of *productions* (aka rewrite rules) of the form:  
non-terminal  $\rightarrow$  (non-terminal | terminal)\*.

Note:

- All sets are finite.
- Left-hand side (lhs) of production consists of a single non-terminal symbol.
- Right-hand side (rhs) consists of arbitrary combinations of terminal and non-terminal symbol described by a regular expression.

# CFG Example and Notation

## Example

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \textit{ident} \end{aligned}$$

## Notation

- Vocabulary  $V = V_t \cup V_n$
- $a, b, c, \dots \in V_t$
- $A, B, C, \dots \in V_n$
- $\alpha, \beta, \gamma, \dots \in V^*$
- $u, v, w, \dots \in V_t^*$



# CFG as a Rewrite System

## Derivations by Rewriting

If  $A \rightarrow \gamma \in P$  then  $\alpha A \beta \Rightarrow \alpha \gamma \beta$  is a *one-step derivation* using  $A \rightarrow \gamma$ .

- We assume that  $\Rightarrow^*$  denotes a derivation of  $\geq 0$  steps.
- In a *leftmost* derivation (denoted  $\Rightarrow_L^*$ ) we replace the leftmost non-terminal in each step.
- In a *rightmost* derivation (denoted  $\Rightarrow_R^*$ ) we replace the rightmost non-terminal in each step.
- If  $S \rightarrow^* \beta$  then  $\beta$  is a *sentential form* of  $G$ .
- If  $S \rightarrow^* w$  and  $w \in V_t^*$  then  $w$  is a *sentence* of  $G$ .
- $L(G) = \{w \mid S \Rightarrow^* w, w \in V_t^*\}$  denotes the language described by  $G$ .

# Example Derivation

## CFG

$$\begin{aligned}E &\rightarrow E + T \mid T \\T &\rightarrow T * F \mid F \\F &\rightarrow (E) \mid \textit{ident}\end{aligned}$$

## Derivation $E \Rightarrow^* \textit{ident} + (\textit{ident})$

$$\begin{aligned}E &\Rightarrow E + T && \Rightarrow T + T \\&\Rightarrow F + T && \Rightarrow \textit{ident} + T \\&\Rightarrow \textit{ident} + F && \Rightarrow \textit{ident} + (E) \\&\Rightarrow \textit{ident} + (T) && \Rightarrow \textit{ident} + (F) \\&\Rightarrow \textit{ident} + (\textit{ident})\end{aligned}$$

# Example Derivation (2)

## CFG

$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid \dots \mid 9$$

## Derivation $E \Rightarrow^* 3 + 7 * 2$

$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E * E \Rightarrow 3 + 7 * E \Rightarrow 3 + 7 * 2$$

## Things to think about

- Multiple derivations for the same input word?
- How to represent derivations?

## Purpose

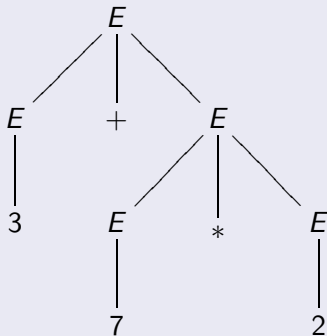
- Compact representation of a derivation  $E \Rightarrow^* w$ .
- Most often in terms of a tree-like structure (but other formats like bit-encodings are also possible).
- Different from AST!
  - Parse tree represents *concrete* input word.
  - AST is an *abstract* representation of input.

# Parse Tree Example

## CFG + Derivation

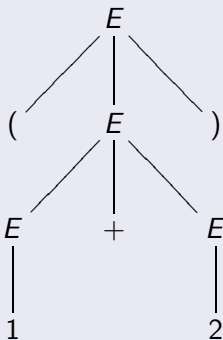
$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid \dots \mid 9$$
$$E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E * E \Rightarrow 3 + 7 * E \Rightarrow 3 + 7 * 2$$

## Parse Tree

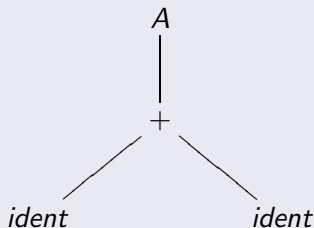


# AST versus Parse Tree

Parse tree



AST



Parse tree: Maintain all source (concrete) syntax info.

AST: Abstract representation.

# Ambiguity

## Definition

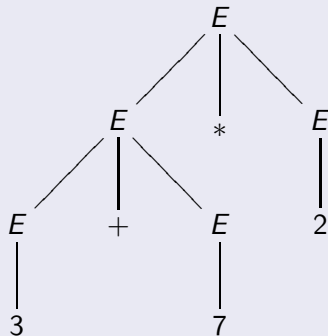
Multiple parse trees for the same grammar and input word.

## Example

$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid \dots \mid 9$$

$$\begin{aligned} E &\Rightarrow E * E \Rightarrow E + E * E \\ &\Rightarrow 3 + E * E \Rightarrow 3 + 7 * E \\ &\Rightarrow 3 + 7 * 2 \end{aligned}$$

Another parse tree:



# Ambiguity Example

## Ambiguous Grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid \text{ident}$$

## Equivalent Non-Ambiguous Grammar

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \text{ident}$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E) \mid \text{ident}$$

## Fact

Some context-free languages are *inherently ambiguous*.



# Ambiguity Example (2)

## Exercise

Consider

$$\begin{array}{lcl} Stmt & \rightarrow & \text{if } Expr \text{ then } Stmt \\ & | & \text{if } Expr \text{ then } Stmt \text{ else } Stmt \\ & | & \text{other} \end{array}$$

Ambiguous? Does an equivalent unambiguous grammar exist?

# Parsing Approaches: Top-Down versus Bottom-Up

## Example

$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid \dots \mid 9$$

## Top-Down (“ANTLR”)

- Build left-most derivation (strictly reduce left-most non-terminal symbols).
- $E \Rightarrow E + E \Rightarrow 3 + E \Rightarrow 3 + E * E \Rightarrow 3 + 7 * E \Rightarrow 3 + 7 * 2$
- “Guess which rule to apply from left to right.”

## Bottom-Up (“yacc”)

- Build right-most derivation.
- $E \Rightarrow E + E \Rightarrow E + E * E \Rightarrow E + E * 2 \Rightarrow E + 7 * 2 \Rightarrow 3 + 7 * 2$
- “Seek for matches of right-hand side, replace by left-hand side.”

# Top-Down Parsing

## Example

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{ident} \end{aligned}$$

## Approach

Build left-most derivation

- starting from start symbol,
- use target string to choose productions.

## Challenge: Left Recursion

- $E \Rightarrow E + T \Rightarrow E + T + T \Rightarrow \dots$
- Several possible right-hand sides.
- Some lead to dead ends. Backtracking? Too inefficient.

# Removal of Left Recursion

## Example

Consider

$$\begin{aligned}E &\rightarrow E + T \mid T \\T &\rightarrow T * F \mid F \\F &\rightarrow (E) \mid \textit{ident}\end{aligned}$$

Remove left recursion, e.g.

$$\begin{aligned}E &\rightarrow TE' \\E' &\rightarrow +TE' \mid \epsilon\end{aligned}$$

# Removal Left Recursion (2)

## Direct Left Recursion Removal

Left recursion can be eliminated as follows:

- Grammar rules

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$

- are replaced by the equivalent rules

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \dots \mid \beta_m A' \\ A' &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \epsilon \end{aligned}$$

## Indirect Left Recursion Removal

Consider

$$A \rightarrow BC$$

$$B \rightarrow A\beta$$

- Left recursion involves several rules.
- Can be eliminated (reduced to direct case) via unfolding.

# Left Factoring

## Issue

Remove common prefix of right hand sides.

## Left Factoring

Left factoring of  $A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \dots \mid \gamma_m$  (where  $\alpha \neq \epsilon$  is the longest common prefix) yields

$$\begin{aligned} A &\rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_m \\ A' &\rightarrow \beta_1 \mid \dots \mid \beta_n \end{aligned}$$

## Example

$$\begin{array}{ll} S \rightarrow iEtS \mid iEtSeS \mid a & \\ E \rightarrow b & \end{array} \quad \text{transformed to} \quad \begin{array}{ll} S \rightarrow iEtSS' \mid a & \\ S' \rightarrow \epsilon \mid eS & \\ E \rightarrow b & \end{array}$$

# Top-Down Parsing: Short Summary

## Predictive Top-Down Parsing

- Eliminate left recursion.
- Apply left factoring.
- Deterministic. No backtracking.

## Lookahead?

- How far to lookahead in the input string to (deterministically) apply a rule?
- Observe right-hand sides!

# First (matching symbols)

## Intuition

Which terminals can start a string matching the non-terminal?

## Definition

$$\text{First}(\alpha) = \{a \in V_t \cup \{\epsilon\} \mid \alpha \Rightarrow^* a\beta\}$$

## Inductive Definition

$$\text{First}(\alpha) = \begin{cases} \{x\} & \text{if } \alpha = x\alpha', x \in V_t \\ \text{First}(x) & \text{if } \alpha = x\alpha', \epsilon \notin \text{First}(x), x \in V_n \\ \text{First}(x) \cup \text{First}(\alpha') & \text{if } \alpha = x\alpha', \epsilon \in \text{First}(x), x \in V_n \end{cases}$$

Sufficient to consider  $\text{First}(X)$  where  $X$  is a non-terminal.



# First Algorithm

## Algorithm

Build  $First(X)$  as follows:

- ❶ If  $X$  is a terminal then  $First(X) = \{X\}$ .
- ❷ If  $X \rightarrow \epsilon$  then add  $\epsilon$  to  $First(X)$ .
- ❸ If  $X \rightarrow Y_1 \dots Y_k$ :
  - ❶ Put  $First(Y_1) - \{\epsilon\}$  in  $First(X)$ .
  - ❷  $\forall i. 1 < i \leq k$ , if  $\epsilon \in First(Y_1) \cap \dots \cap First(Y_{i-1})$  (i.e.  $Y_1 \dots Y_{i-1} \Rightarrow^* \epsilon$ ) then put  $First(Y_i) - \{\epsilon\}$  in  $First(X)$ .
  - ❸ If  $\epsilon \in First(Y_1) \cap \dots \cap First(Y_k)$  then add  $\epsilon$  to  $First(X)$ .

Repeat until no more additions can be made.

# First Example

## Example

$E \rightarrow TE'$	$First(E) = \{ident, (\}$
$E' \rightarrow +TE'$	$First(E') = \{+, \epsilon\}$
$E' \rightarrow \epsilon$	
$T \rightarrow FT'$	$First(T) = \{ident, (\}$
$T' \rightarrow *FT'$	$First(T') = \{*, \epsilon\}$
$T' \rightarrow \epsilon$	
$F \rightarrow ident$	$First(F) = \{ident, (\}$
$F \rightarrow (E)$	

# Follow

## Issue

- Which terminals can start a string matching *after* the nonterminal?
- Consider  $\epsilon$  right hand sides.

## Definition

$$\text{Follow}(A) = \{w \in \text{First}(\beta) \mid S \Rightarrow^* wA\beta\}$$

## Algorithm

Build  $\text{Follow}(A)$  as follows:

- 1 Add  $\$$  to  $\text{Follow}(S)$  where  $S$  is the start symbol.
- 2 If  $A \rightarrow \alpha B \beta$ :
  - 1 Put  $\text{First}(\beta) - \{\epsilon\}$  in  $\text{Follow}(B)$
  - 2 If  $\beta = \epsilon$  (i.e.  $A \rightarrow \alpha B$ ) or  $\epsilon \in \text{First}(\beta)$  (i.e.  $\beta \Rightarrow^* \epsilon$ ) then put  $\text{Follow}(A)$  in  $\text{Follow}(B)$ .

Repeat until no more additions can be made ( $\$$  = EOF token).

# First and Follow Example

## Example

$$S \rightarrow BAa$$

$$A \rightarrow \epsilon \mid BcA$$

$$B \rightarrow \epsilon \mid b$$

	<i>First</i>	<i>Follow</i>
<i>S</i>	<i>a, b, c</i>	<i>\$</i>
<i>A</i>	<i>ϵ, b, c</i>	<i>a</i>
<i>B</i>	<i>ϵ, b</i>	<i>a, b, c</i>

# LL(1) Grammar

## definition

A grammar  $G$  is LL(1) iff for each set of productions  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ :

- 1  $First(\alpha_1), \dots, First(\alpha_n)$  are all pairwise disjoint.
- 2 If  $\alpha_i \Rightarrow^* \epsilon$  then  $First(\alpha_j) \cap Follow(A) = \emptyset$  for all  $1 \leq j \leq n, i \neq j$ .

If  $G$  is  $\epsilon$ -free (there are no  $\epsilon$  productions), first condition is sufficient.

## Facts

- 1 No left-recursive grammar is LL(1).
- 2 No ambiguous grammar is LL(1).
- 3 Some languages have no LL(1) grammar.

# LL(1) Grammar (cont'd)

## Example

$S \rightarrow aS \mid a$  is not LL(1) because  $First(aS) = First(a) = \{a\}$ .  
However, the following equivalent grammar is LL(1).

$$\begin{aligned} S &\rightarrow aS' \\ S' &\rightarrow aS' \mid \epsilon \end{aligned}$$

## Lemma

A grammar is LL(1) iff for each two productions  $A \rightarrow \alpha$  and  $A \rightarrow \beta$  we have that  $Look(A \rightarrow \alpha) \cap Look(A \rightarrow \beta) = \emptyset$  where

$$Look(A \rightarrow \alpha) = \begin{cases} First(\alpha) \setminus \{\epsilon\} & \text{if } \epsilon \notin Follow(\alpha) \\ (First(\alpha) \setminus \{\epsilon\}) \cup Follow(A) & \text{otherwise} \end{cases}$$

# Recursive Descent Parsing

## Example

$E \rightarrow TE'$	$First(E) = \{ident, (\}$
$E' \rightarrow +TE'$	$First(E') = \{+, \epsilon\}$
$E' \rightarrow \epsilon$	
$T \rightarrow FT'$	$First(T) = \{ident, (\}$
$T' \rightarrow *FT'$	$First(T') = \{*, \epsilon\}$
$T' \rightarrow \epsilon$	
$F \rightarrow ident$	$First(F) = \{ident, (\}$
$F \rightarrow (E)$	

## Recursive Descent Parser

Grammar is LL(1). Introduce a function for each non-terminal symbol.

# Recursive Descent Parser

## Parser

Ignore the parse tree, so strictly speaking we only build a matcher.

```
e() { t() ; e'(); if token != EOF then HALT; }
e'() { if token == PLUS
      then get_next_token(); t(); e'(); }
t() { f(); t'(); }
t'() { if token == TIMES
      then get_next_token(); f(); t'(); }
f () { if token == IDENT
      then get_next_token();
      else if token == OPENPAR
          then get_next_token(); e();
          if token==CLOSEPAR
              then get_next_token();
              else HALT; }
```

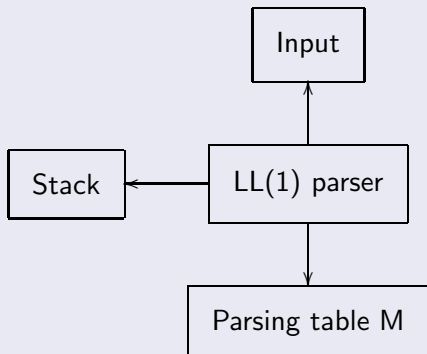


# Table Driven Parsing

## Fact

A program making use of recursive function can be simulated by a non-recursive program (by making use of a stack).

## Table Driven Approach



Stack elements are elements in our vocabulary  
Input consists of terminal symbols  
Table encodes "LL(1) actions"

# Parse Table Construction

## Algorithm

Input: Grammar  $G$

Output: Parsing table  $M$

Method:

- ➊ For each production  $A \rightarrow \alpha$ :
  - ➊ For each  $a \in \text{First}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$ .
  - ➋ If  $\epsilon \in \text{First}(\alpha)$ :
    - ➊ For each  $b \in \text{Follow}(A)$  add  $A \rightarrow \alpha$  to  $M[A, b]$ .
    - ➋ If  $\$ \in \text{Follow}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$ .
- ➋ Set each undefined entry of  $M$  to *error* (or leave empty).

If there are multiple entries then grammar is not LL(1).

# Table Driven Parsing

## Algorithm

```
Push $; Push start symbol
while Top of stack is not $ {
  X:= top of stack; a:= input symbol
  if X is a terminal
    then if X == a
      then pop X; advance input
      else error
    else if M[X,a] == X->Y1 ... Yk
      then pop X; push Yk, ..., Y1
      else error
}
```

# Example

## First and Follow Sets

$E \rightarrow TE'$	$First(E) = \{ident, (\}$	$Follow(E) = \{\$, )\}$
$E' \rightarrow +TE'$	$First(E') = \{+, \epsilon\}$	$Follow(E') = \{\$, )\}$
$E' \rightarrow \epsilon$		
$T \rightarrow FT'$	$First(T) = \{ident, (\}$	$Follow(T) = \{\$, +, )\}$
$T' \rightarrow *FT'$	$First(T') = \{*, \epsilon\}$	$Follow(T') = \{\$, +, )\}$
$T' \rightarrow \epsilon$		
$F \rightarrow ident$	$First(F) = \{ident, (\}$	$Follow(F) = \{\$, +, *, )\}$
$F \rightarrow (E)$		

# Parse Table + Execution

	<i>ident</i>	+	*	(	)	\$
<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>	$F \rightarrow \text{ident}$			$F \rightarrow (E)$		

Stack	Input	Rule
\$E	<i>id</i> + ( <i>id</i> )	$E \rightarrow TE'$
$\$E' T$	<i>id</i> + ( <i>id</i> )	$T \rightarrow FT'$
$\$E' T' F$	<i>id</i> + ( <i>id</i> )	$F \rightarrow \text{ident}$
$\$E' T' \text{id}$	<i>id</i> + ( <i>id</i> )	
$\$E' T'$	+( <i>id</i> )	$T' \rightarrow \epsilon$
$\$E'$	+( <i>id</i> )	$E' \rightarrow +TE'$
$\$E' T+$	+( <i>id</i> )	
$\$E' T$	( <i>id</i> )	$T \rightarrow FT'$
$\$E' T' F$	( <i>id</i> )	$F \rightarrow (E)$
$\$E' T')E($	( <i>id</i> )	

Stack	Input	Rule
$\$E' T')E$	<i>id</i> )	$E \rightarrow TE'$
$\$E' T')E' T$	<i>id</i> )	$T \rightarrow FT'$
$\$E' T')E' T' F$	<i>id</i> )	$F \rightarrow \text{ident}$
$\$E' T')E' T' \text{id}$	<i>id</i> )	
$\$E' T')E' T'$	)	$T' \rightarrow \epsilon$
$\$E' T')E'$	)	$E' \rightarrow \epsilon$
$\$E' T')$	)	
$\$E' T'$	\$	$T' \rightarrow \epsilon$
$\$E'$	\$	$E' \rightarrow \epsilon$
\$	\$	

# Top-Down Parsing

## Summary

- Predictive parsing (removal left recursion, left factoring).
- LL(1) grammars.
- Recursive descent parsing.
- Table driven parsing.
- How to construct parse tree? (later)
- Error recovery important (but no time!), please consult text book.
- ANTRL state-of-the art LL-style parser for Java.