

Part 3

Syntax analysis

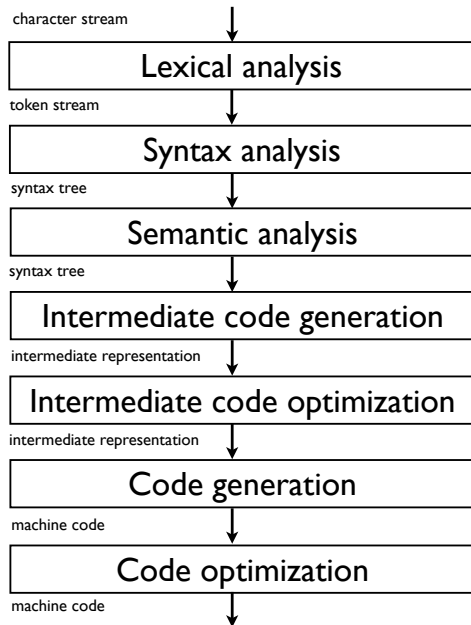
Outline

1. Introduction
2. Context-free grammars
3. Top-down parsing
4. Bottom-up parsing
 - 4.1 Shift/reduce parsing
 - 4.2 LR parsers
 - 4.3 Operator precedence parsing
 - 4.4 Using ambiguous grammars
5. Conclusion and some practical considerations

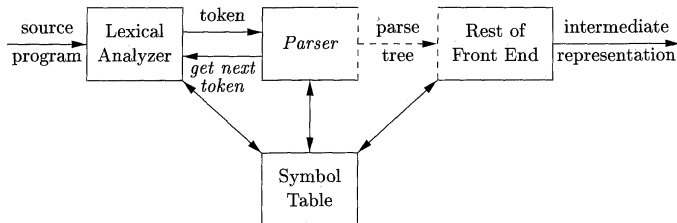
3.1

Introduction

Structure of a compiler



Syntax analysis



■ Goals of the syntax analysis

- ▶ Recombine the tokens provided by the lexical analysis into a structure (called a *syntax tree*)
- ▶ Reject invalid texts by reporting *syntax errors*.

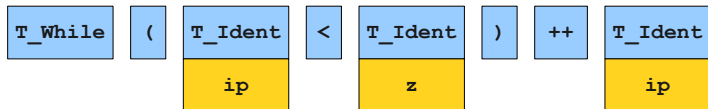
■ Like lexical analysis, syntax analysis is based on

- ▶ the *definition of valid programs* based on some formal languages,
- ▶ the derivation of an *algorithm to detect valid words (programs)* from this language

■ Formal language: context-free grammars

■ Two main algorithm families: Top-down parsing and Bottom-up parsing

Example

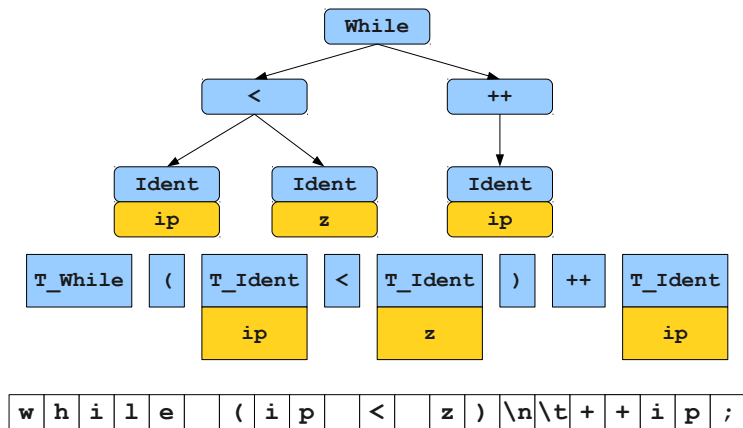


w	h	i	l	e		(i	p		<		z)	\n	\t	+	+	i	p	;
---	---	---	---	---	--	---	---	---	--	---	--	---	---	----	----	---	---	---	---	---

```
while (ip < z)
    ++ip;
```

(Keith Schwarz)

Example



```
while (ip < z)
    ++ip;
```

(Keith Schwarz)

3.2

Context-free grammars

Reminder: grammar

- A **grammar** is a 4-tuple $G = (V, \Sigma, R, S)$, where:
 - ▶ V is an **alphabet**,
 - ▶ $\Sigma \subseteq V$ is the set of **terminal symbols**
($V - \Sigma$ is the set of **nonterminal symbols**),
 - ▶ $R \subseteq (V^+ \times V^*)$ is a finite set of **production rules**
 - ▶ $S \in V - \Sigma$ is the **start symbol**.
- **Notations:**
 - ▶ **Nonterminal symbols** are represented by uppercase letters: A, B, \dots
 - ▶ **Terminal symbols** are represented by lowercase letters: a, b, \dots
 - ▶ **Start symbol** written as S
 - ▶ **Empty word**: ε
 - ▶ A rule $(\alpha, \beta) \in R : \alpha \rightarrow \beta$
 - ▶ **Rule combination**: $A \rightarrow \alpha \mid \beta$

Reminder: grammar

- Example: $\Sigma = \{a, b, c\}$, $V - \Sigma = \{S, R\}$, $R =$

$$S \rightarrow R$$

$$S \rightarrow aSc$$

$$R \rightarrow \varepsilon$$

$$R \rightarrow RbR$$

Reminder: derivation and language

Definitions:

- v can be *derived in one step* from u by G (noted $u \Rightarrow v$) iff $u = xu'y$, $v = xv'y$, and $u' \rightarrow v'$
- v can be *derived in several steps* from u by G (noted $u \Rightarrow^* v$) iff $\exists k \geq 0$ and $v_0 \dots v_k \in V^+$ such that $u = v_0$, $v = v_k$, $v_i \Rightarrow v_{i+1}$ for $0 \leq i < k$
- The *language generated by a grammar* G is the set of words that can be derived from the start symbol:

$$L = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$$

Example: derivation of $aabcc$ from the previous grammar

$$\underline{S} \Rightarrow a\underline{S}c \Rightarrow aa\underline{S}cc \Rightarrow aa\underline{R}cc \Rightarrow aa\underline{R}bRcc \Rightarrow aab\underline{R}cc \Rightarrow aabcc$$

Reminder: type of grammars

Chomsky's grammar hierarchy:

- Type 0: free or unrestricted grammars
- Type 1: context sensitive grammars
 - ▶ productions of the form $\alpha X \beta \rightarrow \alpha \gamma \beta$, where α, β, γ are arbitrary strings of symbols in V , with γ non-null, and X a single nonterminal
- Type 2: context-free grammars (CFG)
 - ▶ productions of the form $X \rightarrow \alpha$ where α is an arbitrary string of symbols in V , and X a single nonterminal.
- Type 3: regular grammars
 - ▶ Productions of the form $X \rightarrow a$, $X \rightarrow aY$ or $X \rightarrow \varepsilon$ where X and Y are nonterminals and a is a terminal (equivalent to regular expressions and finite state automata)

Context-free grammars

- Regular languages are too limited for representing programming languages.
- Examples of languages not representable by a regular expression:
 - ▶ $L = \{a^n b^n \mid n \geq 0\}$
 - ▶ Balanced parentheses
 $L = \{\varepsilon, (), (()), ()(), ((())), (())() \dots\}$
 - ▶ Scheme programs
 $L = \{1, 2, 3, \dots, (\text{lambda}(x)(+ x 1))\}$
- Context-free grammars are typically used for describing programming language syntaxes.
 - ▶ They are sufficient for most languages
 - ▶ They lead to efficient parsing algorithms

Context-free grammars for programming languages

- **Terminals** of the grammars are typically the **tokens derived by the lexical analysis** (in **bold** in rules)
- Divide the language into several **syntactic categories** (sub-languages)
- **Common syntactic categories:**
 - ▶ Expressions: calculation of values
 - ▶ Statements: express actions that occur in a particular sequence
 - ▶ Declarations: express properties of names used in other parts of the program

$Exp \rightarrow Exp + Exp$

$Exp \rightarrow Exp - Exp$

$Exp \rightarrow Exp * Exp$

$Exp \rightarrow Exp / Exp$

$Exp \rightarrow \mathbf{num}$

$Exp \rightarrow \mathbf{id}$

$Exp \rightarrow (Exp)$

$Stat \rightarrow \mathbf{id} := Exp$

$Stat \rightarrow Stat ; Stat$

$Stat \rightarrow \mathbf{if} Exp \mathbf{then} Stat \mathbf{else} Stat$

$Stat \rightarrow \mathbf{if} Exp \mathbf{then} Stat$

Derivation for context-free grammar

- Like for a general grammar
- Because there is only one nonterminal in the LHS of each rule, their order of application does not matter (but the chosen rules do)
- Two particular derivations:
 - ▶ **left-most**: always expand first the left-most nonterminal (important for parsing)
 - ▶ **right-most**: always expand first the right-most nonterminal (canonical derivation)
- Examples:

$$S \rightarrow aTb \mid c$$
$$T \rightarrow cSS \mid S$$
$$w = accacbb$$

Left-most derivation:

$$S \Rightarrow aTb \Rightarrow acSSb \Rightarrow accSb \Rightarrow accaTbb \Rightarrow accaSbb \Rightarrow accacbb$$

Right-most derivation:

$$S \Rightarrow aTb \Rightarrow acSSb \Rightarrow acSaTbb \Rightarrow acSaSbb \Rightarrow acSacbb \Rightarrow accacbb$$

Parse tree

- A **parse tree** abstracts the order of application of the rules
 - ▶ Each **interior node** represents the application of a production
 - ▶ For a rule $A \rightarrow X_1 X_2 \dots X_k$, the interior node is labeled by A and the children from left to right by X_1, X_2, \dots, X_k .
 - ▶ **Leaves** are labeled by nonterminals or terminals and, read from left to right, represent a string generated by the grammar
- A derivation encodes **how** to produce the input
- A parse tree encodes the **structure** of the input
- **Syntax analysis = recovering the parse tree from the tokens**

Parse trees

$$S \rightarrow aTb \mid c$$

$$T \rightarrow cSS \mid S$$

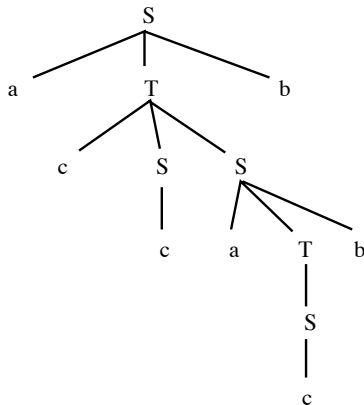
$$w = accacbb$$

Left-most derivation:

$$\begin{aligned} S &\Rightarrow aTb \Rightarrow acSSb \Rightarrow accSb \Rightarrow \\ &accaTbb \Rightarrow accaSbb \Rightarrow accacbb \end{aligned}$$

Right-most derivation:

$$\begin{aligned} S &\Rightarrow aTb \Rightarrow acSSb \Rightarrow acSaTbb \Rightarrow \\ &acSaSbb \Rightarrow acSacbb \Rightarrow accacbb \end{aligned}$$

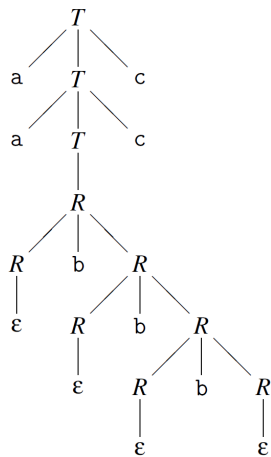
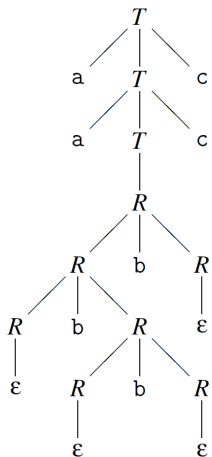


Parse tree

$$T \rightarrow R \mid aTc$$

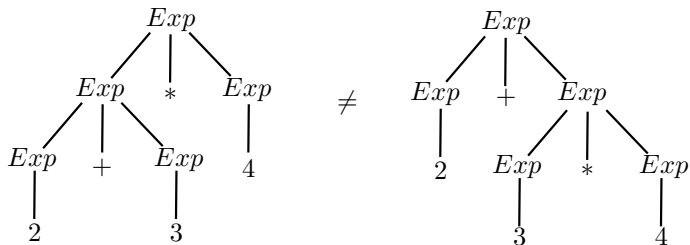
$$R \rightarrow RbR \mid \varepsilon$$

$$w = aabbbcc$$



Ambiguity

- The order of derivation does not matter but the chosen production rules do
- **Definition:** A CFG is **ambiguous** if there is at least one string with **two or more distinct parse trees**
- Ambiguity is not problematic when dealing with flat strings. It is when dealing with language semantics



Detecting and solving Ambiguity

- There is **no mechanical way to determine if a grammar is (un)ambiguous** (this is an undecidable problem)
- In most practical cases however, it is easy to detect and prove ambiguity.
E.g., **any grammar containing $N \rightarrow N\alpha N$ is ambiguous** (two parse trees for $N\alpha N\alpha N$).
- **How to deal with ambiguities?**
 - ▶ **Modify the grammar** to make it unambiguous
 - ▶ **Handle these ambiguities in the parsing algorithm**
- Two common sources of ambiguity in programming languages
 - ▶ Expression syntax (operator precedences)
 - ▶ Dangling else

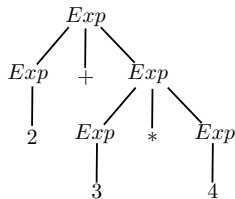
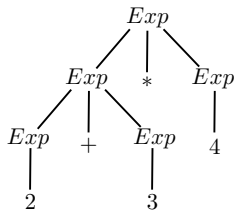
Operator precedence

- This expression grammar is ambiguous

$$Exp \rightarrow Exp + Exp$$
$$Exp \rightarrow Exp - Exp$$
$$Exp \rightarrow Exp * Exp$$
$$Exp \rightarrow Exp / Exp$$
$$Exp \rightarrow \mathbf{num}$$
$$Exp \rightarrow (Exp)$$

(it contains $N \rightarrow N\alpha N$)

- Parsing of $2 + 3 * 4$:



Operator associativity

■ Types of operator associativity:

- ▶ An operator \oplus is **left-associative** if $a \oplus b \oplus c$ must be **evaluated from left to right**, i.e., as $(a \oplus b) \oplus c$
- ▶ An operator \oplus is **right-associative** if $a \oplus b \oplus c$ must be **evaluated from right to left**, i.e., as $a \oplus (b \oplus c)$
- ▶ An operator \oplus is **non-associative** if expressions of the form $a \oplus b \oplus c$ are **not allowed**

■ Examples:

- ▶ $-$ and $/$ are typically **left-associative**
 $(a - b - c = (a - b) - c)$
- ▶ $+$ and $*$ are mathematically **associative (left or right)**.
By convention, we take them left-associative as well
- ▶ **List construction** in functional languages is **right-associative**
- ▶ The **assignment operator** in C is **right-associative**
 $(a = b = c \text{ is equivalent to } a = (b = c))$
- ▶ In Java, **comparison operators** are **non-associative**
(you can not write $2 < 3 < 4$)

Rewriting ambiguous expression grammars

- Let's consider the following ambiguous grammar:

$$E \rightarrow E \oplus E$$

$$E \rightarrow \text{num}$$

- If \oplus is left-associative, we rewrite it as a left-recursive (a recursive reference only to the left).
- If \oplus is right-associative, we rewrite it as a right-recursive (a recursive reference only to the right).

\oplus left-associative

$$E \rightarrow E \oplus E'$$

$$E \rightarrow E'$$

$$E' \rightarrow \text{num}$$

\oplus right-associative

$$E \rightarrow E' \oplus E$$

$$E \rightarrow E'$$

$$E' \rightarrow \text{num}$$

Mixing operators of different precedence levels

- Introduce a different nonterminal for each precedence level

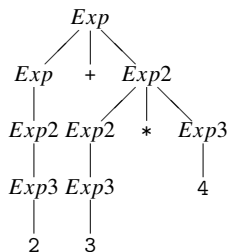
Ambiguous

$Exp \rightarrow Exp + Exp$
 $Exp \rightarrow Exp - Exp$
 $Exp \rightarrow Exp * Exp$
 $Exp \rightarrow Exp / Exp$
 $Exp \rightarrow \text{num}$
 $Exp \rightarrow (Exp)$

Non-ambiguous

$Exp \rightarrow Exp + Term$
 $Exp \rightarrow Exp - Term$
 $Exp \rightarrow Term$
 $Term \rightarrow Term * Factor$
 $Term \rightarrow Term / Factor$
 $Term \rightarrow Factor$
 $Factor \rightarrow \text{num}$
 $Factor \rightarrow (Exp)$

Parse tree for $2 + 3 * 4$



- In this example, the nonterminals Exp , $Term$ and $Factor$ define the application order of the operators in an expression without parentheses.

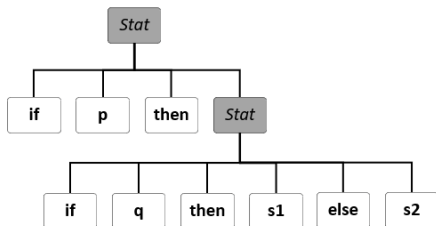
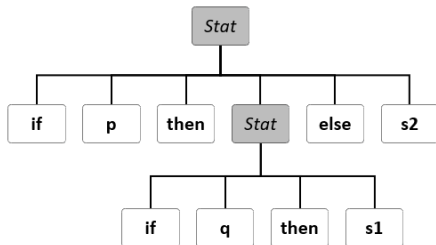
Dangling else

- Else part of a condition is typically optional

Stat → **if** *Exp* **then** *Stat* **else** *Stat*

Stat → **if** *Exp* **then** *Stat*

- How to match `if p then if q then s1 else s2`?



- Convention: `else` matches the closest not previously matched `if`.

Dangling else

■ Unambiguous grammar:

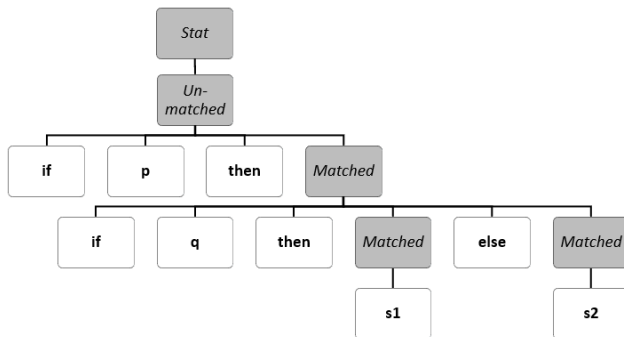
Stat → *Matched* | *Unmatched*

Matched → **if** *Exp* **then** *Matched* **else** *Matched*

Matched → "Any other statement"

Unmatched → **if** *Exp* **then** *Stat*

Unmatched → **if** *Exp* **then** *Matched* **else** *Unmatched*



End-of-file marker

- Parsers must read not only terminal symbols such as $+$, $-$, **num** but also the end-of-file
- We typically use **\$** to represent end of file
- If S is the start symbol of the grammar, then a new start symbol S' is added with the rule $S' \rightarrow S \$$.

$$\begin{aligned} S &\rightarrow Exp \$ \\ Exp &\rightarrow Exp + Term \\ Exp &\rightarrow Exp - Term \\ Exp &\rightarrow Term \\ Term &\rightarrow Term * Factor \\ Term &\rightarrow Term / Factor \\ Term &\rightarrow Factor \\ Factor &\rightarrow \text{num} \\ Factor &\rightarrow (Exp) \end{aligned}$$

Non context-free languages

- Some syntactic constructs from typical programming languages cannot be specified with CFG
- **Example 1:** in C and Java, variables must be declared before their use
 - ▶ $L_1 = \{wcw \mid w \text{ is in } (a \mid b)^*\}$ is not context-free
 - ▶ As a consequence, one can prove that there is no CFG that checks whether a variable is declared before its use
- **Example 2:** checking that a function is called with the right number of arguments
 - ▶ $L_2 = \{a^n b^m c^n d^m \mid n \geq 1 \text{ and } m \geq 1\}$ is not context-free
 - ▶ As a consequence, one can prove that there is no CFG that checks whether a function call is consistent with its definition
 - ▶ In C, the grammar does not count the number of function arguments

$$\begin{aligned} Stmt &\rightarrow \text{ id } (Expr_list) \\ Expr_list &\rightarrow Expr_list, Expr \\ &\mid Expr \end{aligned}$$

- These constructs are typically dealt with during semantic analysis

Backus-Naur Form

- BNF is a text format for describing context-free languages
- We ask you to provide the source grammar for your project in this format
- Example:

```
<expression> ::= <term> | <term> "+" <expression>
<term>        ::= <factor> | <factor> "*" <term>
<factor>      ::= <constant> | <variable> | "(" <expression> ")"
<variable>    ::= "x" | "y" | "z"
<constant>    ::= <digit> | <digit> <constant>
<digit>       ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9"
```

- More information:

http://en.wikipedia.org/wiki/Backus-Naur_form

Syntax analysis

■ Goals:

- ▶ **Checking** that a program is accepted by the context-free grammar
- ▶ Building the **parse tree**
- ▶ **Reporting** syntax errors

■ Two ways:

- ▶ **Top-down**: from the start symbol to the word
- ▶ **Bottom-up**: from the word to the start symbol

Top-down and bottom-up: example

Grammar:

$$S \rightarrow AB$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow b \mid bB$$

Top-down parsing of *aaab*

<i>S</i>	
<i>AB</i>	$S \rightarrow AB$
<i>aAB</i>	$A \rightarrow aA$
<i>aaAB</i>	$A \rightarrow aA$
<i>aaaAB</i>	$A \rightarrow aA$
<i>aaaεB</i>	$A \rightarrow \varepsilon$
<i>aaab</i>	$B \rightarrow b$

Bottom-up parsing of *aaab*

<i>aaab</i>	
<i>aaaεb</i>	(insert ε)
<i>aaaAb</i>	$A \rightarrow \varepsilon$
<i>aaAb</i>	$A \rightarrow aA$
<i>aAb</i>	$A \rightarrow aA$
<i>Ab</i>	$A \rightarrow aA$
<i>AB</i>	$B \rightarrow b$
<i>S</i>	$S \rightarrow AB$

3.3

Top-down parsing

A naive top-down parser

■ A very naive parsing algorithm:

- ▶ Generate **all possible parse trees** until you get one that matches your input
- ▶ To generate all parse trees:
 1. Start with the root of the parse tree (the start symbol of the grammar)
 2. Choose a non-terminal A at one leaf of the current parse tree
 3. Choose a production having that non-terminal as LHS, eg.,
 $A \rightarrow X_1 X_2 \dots X_k$
 4. Expand the tree by making X_1, X_2, \dots, X_k , the children of A .
 5. Repeat at step 2 until all leaves are terminals
 6. Repeat the whole procedure by changing the productions chosen at step 3

(Note: the choice of the non-terminal in Step 2 is irrelevant for a context-free grammar)

- This algorithm is **very inefficient**, **does not always terminate**, etc.

Top-down parsing with backtracking

- Modifications of the previous algorithm:
 1. Depth-first development of the parse tree (corresponding to a left-most derivation)
 2. Process the terminals in the RHS during the development of the tree, checking that they match the input
 3. If they don't at some step, stop expansion and restart at the previous non-terminal with another production rules (**backtracking**)
- Depth-first can be implemented by storing the unprocessed symbols on a **stack**
- Because of the left-most derivation, the inputs can be processed from left to right

Backtracking example

	Stack	Inputs	Action
	<i>S</i>	<i>bcd</i>	Try $S \rightarrow bab$
	<i>bab</i>	<i>bcd</i>	match <i>b</i>
$S \rightarrow bab$	<i>ab</i>	<i>cd</i>	dead-end, backtrack
$S \rightarrow bA$	<i>S</i>	<i>bcd</i>	Try $S \rightarrow bA$
$A \rightarrow d$	<i>bA</i>	<i>bcd</i>	match <i>b</i>
$A \rightarrow cA$	<i>A</i>	<i>cd</i>	Try $A \rightarrow d$
	<i>d</i>	<i>cd</i>	dead-end, backtrack
	<i>A</i>	<i>cd</i>	Try $A \rightarrow cA$
	<i>cA</i>	<i>cd</i>	match <i>c</i>
$w = bcd$	<i>A</i>	<i>d</i>	Try $A \rightarrow d$
	<i>d</i>	<i>d</i>	match <i>d</i>
			Success!

Top-down parsing with backtracking

- General algorithm (to match a word w):

Create a stack with the start symbol

$X = \text{POP}()$

$a = \text{GETNEXTTOKEN}()$

while (True)

if (X is a nonterminal)

 Pick next rule to expand $X \rightarrow Y_1 Y_2 \dots Y_k$

 Push Y_k, Y_{k-1}, \dots, Y_1 on the stack

$X = \text{POP}()$

elseif ($X == \$$ and $a == \$$)

 Accept the input

elseif ($X == a$)

$a = \text{GETNEXTTOKEN}()$

$X = \text{POP}()$

else

 Backtrack

- Ok for small grammars but **still untractable and very slow for large grammars**
- **Worst-case exponential time** in case of syntax error

Top-down parsing with backtracking

- A top-down parser (even with backtracking) can go into an **infinite loop** if the grammar is **left-recursive**:

$$Exp \rightarrow Exp + Term \mid Term$$

$$Term \rightarrow Term * Factor \mid Factor$$

$$Factor \rightarrow num \mid id \mid (Exp)$$

$$w = (id + id) * num$$

Stack	Inputs	Action
<i>Exp</i>	$(id + id) * num$	Try $Exp \rightarrow Exp + Term$
$Exp + Term$	$(id + id) * num$	Try $Exp \rightarrow Exp + Term$
$Exp + Term + Term$	$(id + id) * num$	Try $Exp \rightarrow Exp + Term$
...

- Such grammars have to be **reformulated** in order to avoid left-recursion (see later)

Another example

$S \rightarrow aSbT$

$S \rightarrow cT$

$S \rightarrow d$

$T \rightarrow aT$

$T \rightarrow bS$

$T \rightarrow c$

$w = accbbadbc$

Stack	Inputs	Action
S	$accbbadbc$	Try $S \rightarrow aSbT$
$aSbT$	$accbbadbc$	match a
SbT	$accbbadbc$	Try $S \rightarrow aSbT$
$aSbTbT$	$accbbadbc$	match a
$SbTbT$	$ccbbadbc$	Try $S \rightarrow cT$
$cTbTbT$	$ccbbadbc$	match c
$TbTbT$	$cbbadbc$	Try $T \rightarrow c$
$cbTbT$	$cbbadbc$	match cb
TbT	$badbc$	Try $T \rightarrow bS$
$bSbT$	$badbc$	match b
SbT	$adbc$	Try $S \rightarrow aSbT$
$aSbT$	$adbc$	match a
...
c	c	match c
		Success!

Predictive parsing

■ Predictive parser:

- ▶ In the previous example, the production rule to apply can be **predicted** based solely **on the next input symbol and the current nonterminal**
- ▶ **Much faster** than backtracking but this trick works **only for some specific grammars**

■ Grammars for which top-down predictive parsing is possible by looking at the next symbol are called **LL(1) grammars**:

- ▶ L: left-to-right scan of the tokens
- ▶ L: leftmost derivation
- ▶ (1): one token of lookahead

■ Predicted rules are stored in a **parsing table** M :

- ▶ $M[X, a]$ stores the rule to apply when the nonterminal X is on the stack and the next input terminal is a

Example: parse table

$S \rightarrow E\$$

$E \rightarrow \text{int}$

$E \rightarrow (E \text{ Op } E)$

$\text{Op} \rightarrow +$

$\text{Op} \rightarrow *$

	int	()	+	*	\$
S	E\$	E\$				
E	int	(E Op E)				
Op				+	*	

(Keith Schwarz)

Example: successful parsing

1. $S \rightarrow E\$$
2. $E \rightarrow \text{int}$
3. $E \rightarrow (E \text{ Op } E)$
4. $\text{Op} \rightarrow +$
5. $\text{Op} \rightarrow -$

	int	()	+	*	\$
S	1	1				
E	2	3				
Op				4	5	

S	(int + (int * int))\$
E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$
int))\$	int))\$
)\$)\$
)\$)\$
\$	\$

(Keith Schwarz)

Example: erroneous parsing

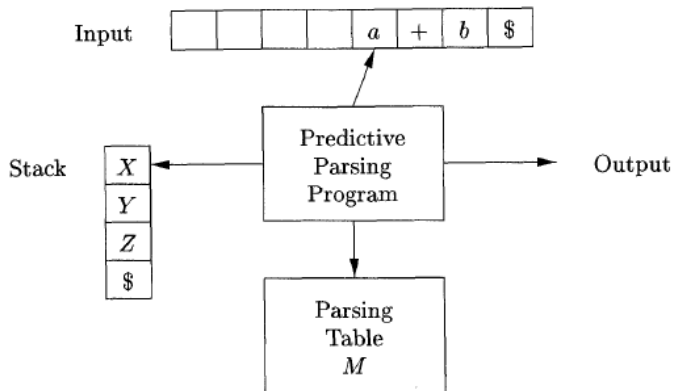
1. $S \rightarrow E\$$
2. $E \rightarrow \text{int}$
3. $E \rightarrow (E \text{ Op } E)$
4. $\text{Op} \rightarrow +$
5. $\text{Op} \rightarrow -$

S	(int (int))\$
E\$	(int (int))\$
(E Op E) \$	(int (int))\$
E Op E) \$	int (int))\$
int Op E) \$	int (int))\$
Op E) \$	(int))\$

	int	()	+	*	\$
S	1	1				
E	2	3				
Op				4	5	

(Keith Schwarz)

Table-driven predictive parser



(Dragonbook)

Table-driven predictive parser

Create a stack with the start symbol

$X = \text{POP}()$

$a = \text{GETNEXTTOKEN}()$

while (True)

if (X is a nonterminal)

if ($M[X, a] == \text{NULL}$)

 Error

elseif ($M[X, a] == X \rightarrow Y_1 Y_2 \dots Y_k$)

 Push Y_k, Y_{k-1}, \dots, Y_1 on the stack

$X = \text{POP}()$

elseif ($X == \$$ and $a == \$$)

 Accept the input

elseif ($X == a$)

$a = \text{GETNEXTTOKEN}()$

$X = \text{POP}()$

else

 Error

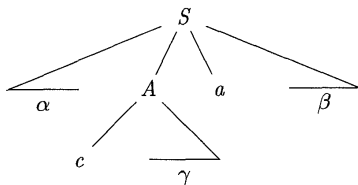
$LL(1)$ grammars and parsing

Three questions we need to address:

- How to build the table for a given grammar?
- How to know if a grammar is $LL(1)$?
- How to change a grammar to make it $LL(1)$?

Building the table

- It is useful to define three functions
(with A a nonterminal and α any sequence of grammar symbols):
 - ▶ **Nullable**(α) is true if $\alpha \Rightarrow^* \varepsilon$
 - ▶ **First**(α) returns the set of terminals c such that $\alpha \Rightarrow^* c\gamma$ for some (possibly empty) sequence γ of grammar symbols
 - ▶ **Follow**(A) returns the set of terminals a such that $S \Rightarrow^* \alpha A a \beta$, where α and β are (possibly empty) sequences of grammar symbols



($c \in \text{First}(A)$ and $a \in \text{Follow}(A)$)

Example

$C \rightarrow P F \text{ class id } X Y \$$

$P \rightarrow \text{public} \mid \varepsilon$

$F \rightarrow \text{final} \mid \varepsilon$

$X \rightarrow \text{extends id} \mid \varepsilon$

$Y \rightarrow \text{implements } I \mid \varepsilon$

$I \rightarrow \text{id } J$

$J \rightarrow , I \mid \varepsilon$

$First(C) = \{ \text{public} , \text{final} , \text{class} \}$

$First(P) = \{ \text{public} \}$

$First(F) = \{ \text{final} \}$

$First(X) = \{ \text{extends} \}$

$First(I) = \{ \text{id} \}$

$First(J) = \{ ' , ' \}$

$Follow(C) = \{ \}$

$Follow(P) = \{ \text{final} , \text{class} \}$

$Follow(F) = \{ \text{class} \}$

$Follow(X) = \{ \text{implements} , \$ \}$

$Follow(Y) = \{ \$ \}$

$Follow(I) = \{ \$ \}$

$Follow(J) = \{ \$ \}$

(Carnahan)

Building the table from *First*, *Follow*, and *Nullable*

To construct the table:

1. Start with the empty table
2. For each production $A \rightarrow \alpha$:
 - 2.1 add $A \rightarrow \alpha$ to $M[A, a]$ for each terminal a in $First(\alpha)$
 - 2.2 if $Nullable(\alpha)$ then add $A \rightarrow \alpha$ to $M[A, a]$ for each a in $Follow(A)$

Illustration:

	public	final	class	...
C	P F class id X Y \$	P F class id X Y \$	P F class id X Y \$	
P	public	ϵ	ϵ	
F		final	ϵ	
...				

LL(1) grammars

■ Three situations:

- ▶ $M[A, a]$ is **empty**: no production is appropriate. We can not parse the sentence and have to report a **syntax error**
- ▶ $M[A, a]$ contains **one entry**: **perfect!**
- ▶ $M[A, a]$ contains **two entries**: **the grammar is not appropriate** for predictive parsing (with one token lookahead)

■ **Definition:** A grammar is **LL(1)** if its parsing table contains at most one entry in each cell or, equivalently, if for all production pairs $A \rightarrow \alpha \mid \beta$

- ▶ $First(\alpha) \cap First(\beta) = \emptyset$,
- ▶ $Nullable(\alpha)$ and $Nullable(\beta)$ are not both true,
- ▶ if $Nullable(\beta)$, then $First(\alpha) \cap Follow(A) = \emptyset$

Example of a non $LL(1)$ grammar

$Exp \rightarrow Exp + Term$

$Exp \rightarrow Exp - Term$

$Exp \rightarrow Term$

$Term \rightarrow Term * Factor$

$Term \rightarrow Term / Factor$

$Term \rightarrow Factor$

$Factor \rightarrow \mathbf{num}$

$Factor \rightarrow (Exp)$

($First(\alpha)$ is the same for all RHS α of the productions for Exp and $Term$)

Computing *Nullable*

Algorithm to compute *Nullable* for all grammar symbols:

Initialize *Nullable* to *False*.

repeat

for each production $X \rightarrow Y_1 Y_2 \dots Y_k$

if $Y_1 \dots Y_k$ are all nullable (or if $k = 0$)

$Nullable(X) = True$

until *Nullable* did not change in this iteration.

Algorithm to compute *Nullable* for any string $\alpha = X_1 X_2 \dots X_k$:

if $X_1 \dots X_k$ are all nullable

$Nullable(\alpha) = True$

else

$Nullable(\alpha) = False$

Computing *First*

Algorithm to compute *First* for all grammar symbols:

Initialize *First* to empty sets.

for each terminal z

$$First(z) = \{z\}$$

repeat

for each production $X \rightarrow Y_1 Y_2 \dots Y_k$

for $i = 1$ **to** k

if $Y_1 \dots Y_{i-1}$ are all nullable (or $i = 1$)

$$First(X) = First(X) \cup First(Y_i)$$

until *First* did not change in this iteration.

Algorithm to compute *First* for any string $\alpha = X_1 X_2 \dots X_k$:

Initialize $First(\alpha) = \emptyset$

for $i = 1$ **to** k

if $X_1 \dots X_{i-1}$ are all nullable (or $i = 1$)

$$First(\alpha) = First(\alpha) \cup First(X_i)$$

Computing *Follow*

Algorithm to compute *Follow* for all nonterminal symbols:

Initialize *Follow* to empty sets.

repeat

for each production $X \rightarrow Y_1 Y_2 \dots Y_k$

for $i = 1$ **to** k , **for** $j = i + 1$ **to** k

if $Y_{i+1} \dots Y_k$ are all nullable (or $i = k$)

$Follow(Y_i) = Follow(Y_i) \cup Follow(X)$

if $Y_{i+1} \dots Y_{j-1}$ are all nullable (or $i + 1 = j$)

$Follow(Y_i) = Follow(Y_i) \cup First(Y_j)$

until *Follow* did not change in this iteration.

Example

Compute the parsing table for the following grammar:

$$S \rightarrow E\$$$

$$E \rightarrow T E'$$

$$E' \rightarrow + T E'$$

$$E' \rightarrow - T E'$$

$$E' \rightarrow \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T'$$

$$T' \rightarrow / F T'$$

$$T' \rightarrow \varepsilon$$

$$F \rightarrow \mathbf{id}$$

$$F \rightarrow \mathbf{num}$$

$$F \rightarrow (E)$$

Example

Nonterminals	Nullable	First	Follow
S	False	{(, id , num }	\emptyset
E	False	{(, id , num }	{), \$}
E'	True	{+, -}	{), \$}
T	False	{(, id , num }	{), +, -, \$}
T'	True	{*, /}	{), +, -, \$}
F	False	{(, id , num }	{), *, /, +, -, \$}

	+	*	id	()	\$
S			$S \rightarrow E\$$	$S \rightarrow E\$$		
E			$E \rightarrow TE'$	$E \rightarrow TE'$		
E'	$E' \rightarrow +TE'$				$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T			$T \rightarrow FT'$	$T \rightarrow FT'$		
T'	$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$			$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F			$F \rightarrow \mathbf{id}$	$F \rightarrow (E)$		

(-, /, and **num** are treated similarly)

Transforming a grammar for $LL(1)$ parsing

- Ambiguous grammars are not $LL(1)$ but **unambiguous grammars are not necessarily $LL(1)$**
- Having a non- $LL(1)$ unambiguous grammar for a language does not mean that this language is not $LL(1)$.
- But there are languages for which there exist unambiguous context-free grammars but no $LL(1)$ grammar.
- We will see two **grammar transformations that improve the chance to get a $LL(1)$ grammar**:
 - ▶ Elimination of left-recursion
 - ▶ Left-factorization

Left-recursion

- The following expression grammar is unambiguous but it is not $LL(1)$:

$$Exp \rightarrow Exp + Term$$
$$Exp \rightarrow Exp - Term$$
$$Exp \rightarrow Term$$
$$Term \rightarrow Term * Factor$$
$$Term \rightarrow Term / Factor$$
$$Term \rightarrow Factor$$
$$Factor \rightarrow \mathbf{num}$$
$$Factor \rightarrow (Exp)$$

- Indeed, $First(\alpha)$ is the same for all RHS α of the productions for Exp et $Term$
- This is a consequence of *left-recursion*.

Left-recursion

- **Recursive productions** are productions defined in terms of themselves. Examples: $A \rightarrow Ab$ or $A \rightarrow bA$.
- When the recursive nonterminal is at the left (resp. right), the production is said to be **left-recursive** (resp. **right-recursive**).
- **Left-recursive productions can be rewritten with right-recursive productions**
- Example:

$$\begin{array}{l} N \rightarrow N\alpha_1 \\ \vdots \\ N \rightarrow N\alpha_m \\ N \rightarrow \beta_1 \\ \vdots \\ N \rightarrow \beta_n \end{array} \quad \Leftrightarrow \quad \begin{array}{l} N \rightarrow \beta_1 N' \\ \vdots \\ N \rightarrow \beta_n N' \\ N' \rightarrow \alpha_1 N' \\ \vdots \\ N' \rightarrow \alpha_m N' \\ N' \rightarrow \varepsilon \end{array}$$

Right-recursive expression grammar

$$Exp \rightarrow Exp + Term$$
$$Exp \rightarrow Exp - Term$$
$$Exp \rightarrow Term$$
$$Term \rightarrow Term * Factor$$
$$Term \rightarrow Term / Factor$$
$$Term \rightarrow Factor$$
$$Factor \rightarrow \mathbf{num}$$
$$Factor \rightarrow (Exp)$$
$$Exp \rightarrow Term Exp'$$
$$Exp' \rightarrow + Term Exp'$$
$$Exp' \rightarrow - Term Exp'$$
$$Exp' \rightarrow \varepsilon$$
$$Term \rightarrow Factor Term'$$
$$Term' \rightarrow * Factor Term'$$
$$Term' \rightarrow / Factor Term'$$
$$Term' \rightarrow \varepsilon$$
$$\Leftrightarrow$$
$$Factor \rightarrow \mathbf{num}$$
$$Factor \rightarrow (Exp)$$

Left-factorisation

- The RHS of these two productions have the same *First* set.

$$\text{Stat} \rightarrow \text{if Exp then Stat else Stat}$$
$$\text{Stat} \rightarrow \text{if Exp then Stat}$$

- The problem can be solved by **left factorising** the grammar:

$$\text{Stat} \rightarrow \text{if Exp then Stat ElseStat}$$
$$\text{ElseStat} \rightarrow \text{else Stat}$$
$$\text{ElseStat} \rightarrow \varepsilon$$

- **Note**

- ▶ The resulting grammar is ambiguous and the parsing table will contain two rules for $M[\text{ElseStat}, \text{else}]$ (because $\text{else} \in \text{Follow}(\text{ElseStat})$ and $\text{else} \in \text{First}(\text{else Stat})$)
- ▶ Ambiguity can be solved in this case by letting $M[\text{ElseStat}, \text{else}] = \{\text{ElseStat} \rightarrow \text{else Stat}\}.$

Hidden left-factors and hidden left recursion

- Sometimes, left-factors or left recursion are hidden

- Examples:

- ▶ The following grammar:

$$A \rightarrow da \mid acB$$

$$B \rightarrow abB \mid daA \mid Af$$

has two overlapping productions: $B \rightarrow daA$ and $B \xRightarrow{*} daf$.

- ▶ The following grammar:

$$S \rightarrow Tu \mid wx$$

$$T \rightarrow Sq \mid vvS$$

has left recursion on T ($T \xRightarrow{*} Tuq$)

- Solution: expand the production rules by substitution to make left-recursion or left factors visible and then eliminate them

Summary

Construction of a $LL(1)$ parser from a CFG grammar

- Eliminate ambiguity
- Eliminate left recursion
- Left factorization
- Add an extra start production $S' \rightarrow S\$$ to the grammar
- Calculate *First* for every production and *Follow* for every nonterminal
- Calculate the parsing table
- Check that the grammar is $LL(1)$

Recursive implementation

- From the parsing table, it is easy to implement a predictive parser recursively (with one function per nonterminal)

$$\begin{aligned}T' &\rightarrow T\$ \\T &\rightarrow R \\T &\rightarrow aTc \\R &\rightarrow \epsilon \\R &\rightarrow bR\end{aligned}$$

	a	b	c	\$
T'	$T' \rightarrow T\$$	$T' \rightarrow T\$$		$T' \rightarrow T\$$
T	$T \rightarrow aTc$	$T \rightarrow R$	$T \rightarrow R$	$T \rightarrow R$
R		$R \rightarrow bR$	$R \rightarrow \epsilon$	$R \rightarrow \epsilon$

```
function parseT'() =  
  if next = 'a' or next = 'b' or next = '$' then  
    parseT() ; match('$')  
  else reportError()
```

```
function parseT() =  
  if next = 'b' or next = 'c' or next = '$' then  
    parseR()  
  else if next = 'a' then  
    match('a') ; parseT() ; match('c')  
  else reportError()
```

```
function parseR() =  
  if next = 'c' or next = '$' then  
    (* do nothing *)  
  else if next = 'b' then  
    match('b') ; parseR()  
  else reportError()
```

(Mogensen)