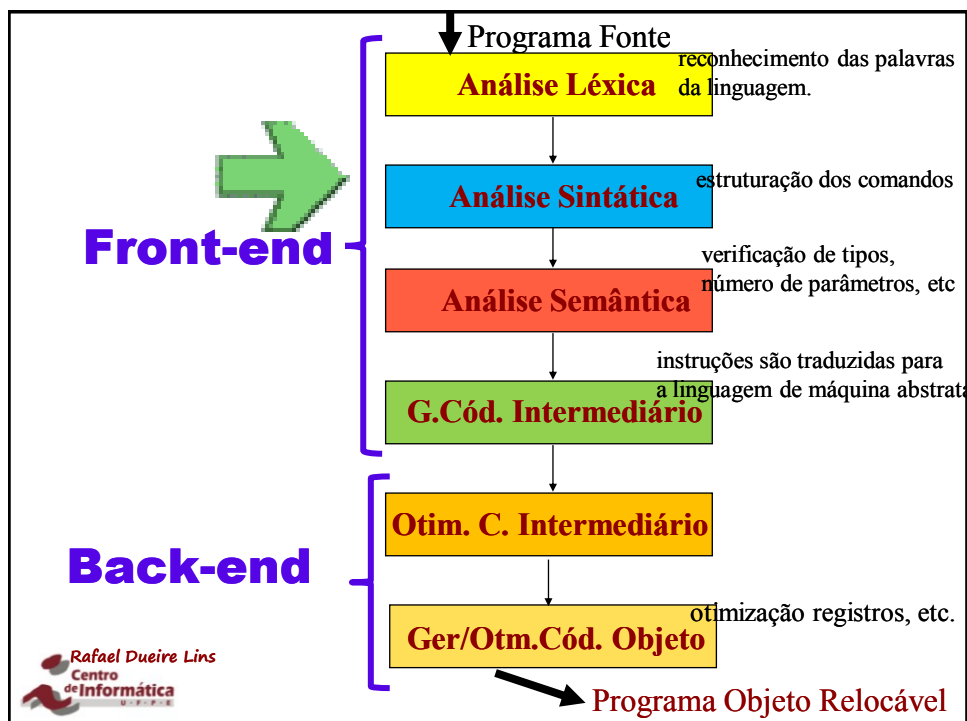
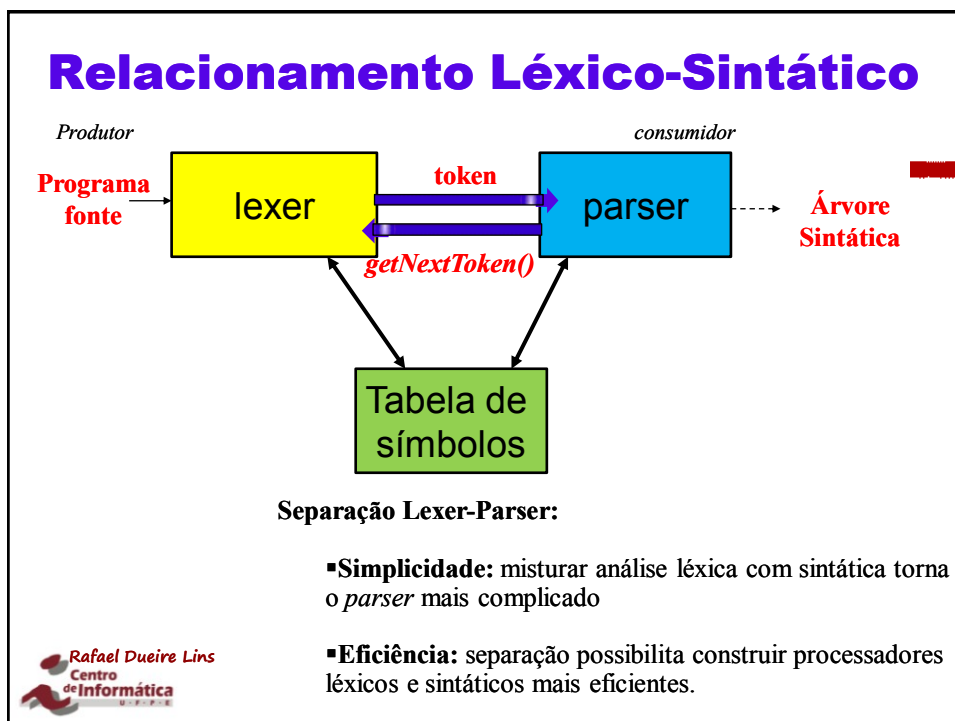
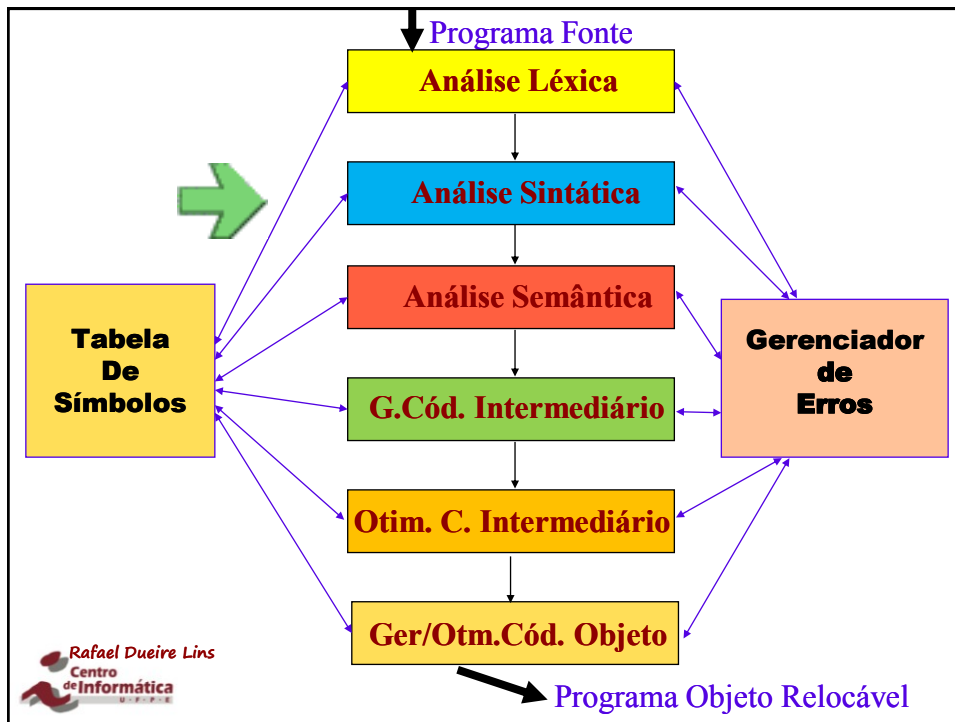


# Análise Sintática

**Rafael Dueire Lins**

*A análise sintática ou parsing é a segunda  
Fase da compilação.*





## Gramática

### Sintaxe de uma linguagem

- Define os *strings* estruturalmente válidos de uma linguagem.
- *Parsing* correto:  
Todos as sentenças estão “bem-formadas”.
- O programa pode ainda...
  - Não passar na checagem de tipos (<70 % erros)
  - Conter erros na lógica!!!



## Gramática

### Especificação da Sintaxe

- Dever ser precisa e fácil de entender
  - Precisa: não ambígua
  - Fácil de entender: deixa evidente a sintaxe da linguagem
- Gramáticas livre de contexto.
- Forma de Backus-Naur (BNFs) é um formalismo adequado
  - Fácil de especificar/manter/entender



# Gramática

## Especificação da Sintaxe

- **Terminais:**
  - Símbolos básicos com os quais *strings* são formados.
  - Token  $\equiv$  terminal (if, then, else, while, etc.)
- **Não-terminais:**
  - Variáveis sintáticas – conjuntos de strings.
  - Símbolo inicial.
- **Produções:**
  - Modo como terminais e não-terminais formam *strings*.



# Gramática

## Convenções Notacionais

- **Terminais:**
  - a, b, c, ... (minúsculas iniciais do alfabeto)
  - +, -, \*, (operadores)
  - (, ), “, ”, ; (símbolos de pontuação)
  - 0, 1, 2, ..., 9 (dígitos)
  - Strings em negrito (e.g.: **id**, **if**, **then**)
- **Não-terminais:**
  - A, B, C, ... (maiúsculas iniciais do alfabeto)
  - S – geralmente o símbolo inicial.
  - Nomes minúsculos em itálico (*expr*, *stmt*, ...)
- **Símbolos da Gramática:**
  - X, Y, Z - terminais ou não-terminais



# Gramática

## Convenções Notacionais

- **Strings de Terminais:**  
u, v, ..., z, ... (minúsculas finais do alfabeto)
- **Strings de símbolos da gramática:**  
 $\alpha, \beta, \gamma, \delta, \eta, \dots$  (minúsculas do alfabeto grego)
- **Alternativas:**  
 $A \rightarrow \alpha_1, A \rightarrow \alpha_2, A \rightarrow \alpha_3, A \rightarrow \alpha_n,$   
 $A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3 \mid \dots \mid \alpha_n,$
- **Símbolo inicial:**  
Se nada dito em contrário, o lado esquerdo da primeira produção.

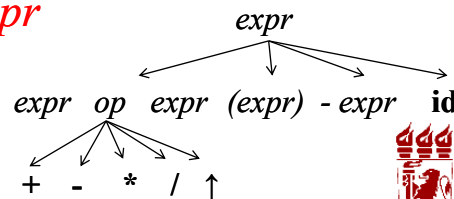


# Gramática

## Exemplo

$expr \rightarrow expr \ op \ expr \mid (expr) \mid - \ expr \mid id$   
 $op \rightarrow + \mid - \mid * \mid / \mid \uparrow$

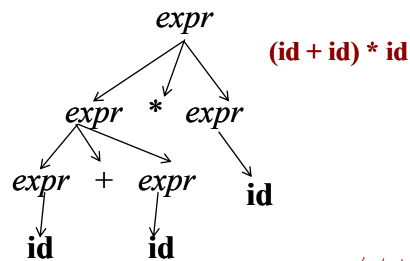
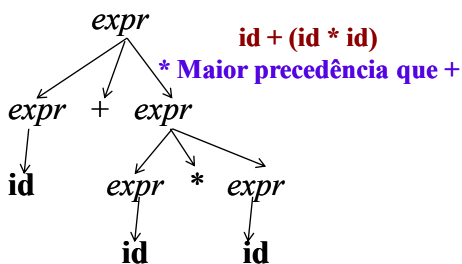
- **Terminais:**  
+ - \* / id  $\uparrow$
- **Não-terminais:** *expr, op.*
- **Símbolo inicial:** *expr*



## Gramática Ambígua Gera Mais de Uma Árvore Sintática

$expr \rightarrow expr\ op\ expr \mid (expr) \mid -\ expr \mid id$   
 $op \rightarrow + \mid - \mid * \mid / \mid \uparrow$

**id + id \* id**



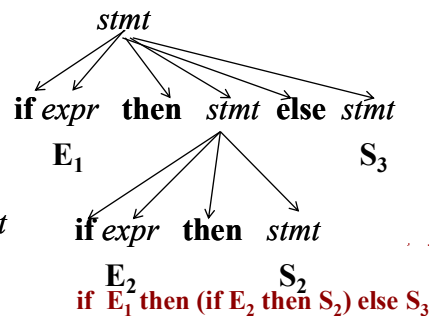
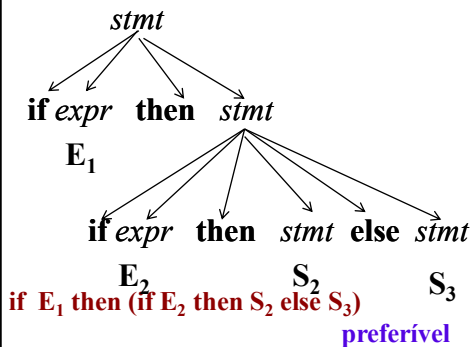
Ambiguidade deve ser evitada!!!



## Gramática Ambígua Gera Mais de uma Árvore Sintática

$stmt \rightarrow \text{if } expr \text{ then } stmt$   
 $\quad \mid \text{if } expr \text{ then } stmt \text{ else } stmt$   
 $\quad \mid \text{other}$

**if  $E_1$  then if  $S_1$  then  $S_2$  else  $S_3$**



## Recursão à esquerda

- Dificuldade em saber quando parar de aplicar uma produção

$S \rightarrow aABe$        $a\underline{A}Be$   
 $A \rightarrow A\underline{b}c \mid b$        $\rightarrow a\underline{A}bcbBe$   
 $B \rightarrow d$        $\rightarrow a\underline{A}bcbcbBe$   
       $\rightarrow a\underline{A}bcbcbcbBe$   
       $\rightarrow a\underline{A}bcbcbcbcbBe$   
       $\rightarrow \dots$



## Expressões Regulares vs. Gramáticas Livres de Contexto

- Toda ER pode ser descrita como uma GLC.
- Conversão automática  $ER \rightarrow GLC$ .
- Por que usar ER para a parte léxica?
  - Regras léxicas são muito simples!
  - ER são mais concisas e intuitivas.
  - ER permitem analisadores ótimos!
  - Modularidade e Mantainabilidade!
- Estruturas aninhadas não podem ser expressas por ER (não há memória).

## Construtores Não-Livres-de-Contexto

- Algumas linguagens não podem ser especificadas por gramáticas somente.
- $L_1 = \{wcw \mid w \text{ em } (a|b)^*\}$   
declaração – uso do identificador
- $L_2 = \{a^n b^m c^n d^m \mid n \geq 1 \text{ e } m \geq 1\}$   
contagem dos parâmetros em uma chamada de subrotina
- $L_3 = \{a^n b^n c^n \mid n \geq 0\}$   
contagem dos parâmetros
- Em alguns casos pode-se obter GLCs equivalentes!!

## Tipos de *Parsers*

- Universais:
  - todo tipo de gramática.
  - Cocke-Younger[1967]-Kasami [1965] e Earley [1970]
  - Ineficientes
- Top-down ou Bottom-up:
  - Sub-classes de gramáticas tais como LL e LR.
  - LR podem ser automatizados.



## Top-down parsing

- Da raiz para as folhas.
- Procura sequência de derivações mais à esquerda para se obter um *string* de entrada
- O parsing do *string* **abbcbcbde** é feito pela sequência de derivações abaixo:

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

$S \rightarrow a\underline{A}Be$   
 $\rightarrow a\underline{A}bcBe$   
 $\rightarrow a\underline{A}bcbcbBe$   
 $\rightarrow abbcbcb\underline{B}e$   
 $\rightarrow abbcbcbde$



## Top-down parsing

- Inicie com símbolo raiz
- Consuma *tokens* da esquerda para à direita
- Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e S



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e

aABe



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e

AB  
e

escolha!



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid \mathbf{b}$

$B \rightarrow d$

a b b c b c d e

bBe



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid \mathbf{b}$

$B \rightarrow d$

a b b c b c d e

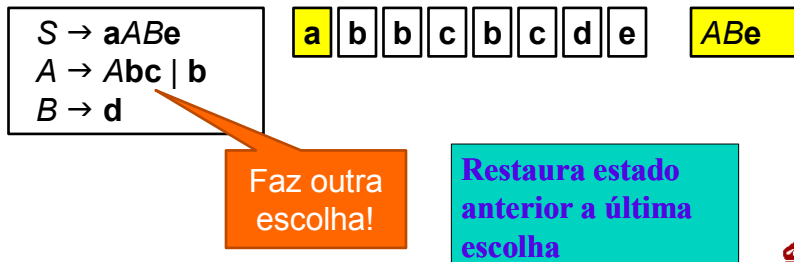
Be

✗ Erro. Backtrack!



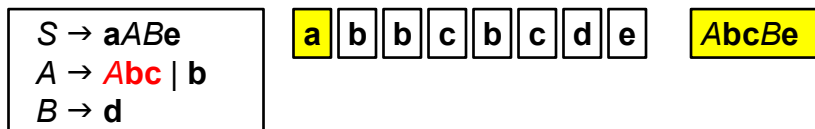
## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar



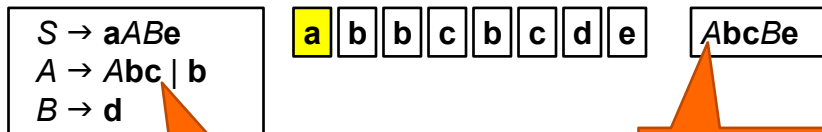
## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

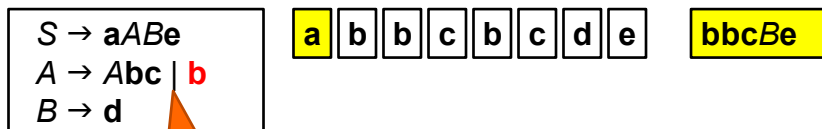


escolha novamente!

Símbolo A novamente na posição mais à esquerda

## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar



Outra escolha!

## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

a b b c b c d e

bcBe



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

a b b c b c d e

Be

✗ Erro. Backtrack!



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e

AbcBe

Restaura estado  
anterior a última  
escolha



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e

AbcbcbBe

Faz outra  
escolha!



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

a b b c b c d e

bbcbbcBe

## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

a b b c b c d e

Be



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e de



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e



## Top-down parsing

- ⌘ Inicie com símbolo raiz
- ⌘ Consuma tokens da esquerda para direita
- ⌘ Decida que produção aplicar

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

a b b c b c d e

derivação  
correspondente!

$S \rightarrow a\text{ABe$   
 $\rightarrow a\text{AbcBe}$   
 $\rightarrow a\text{AbcbcbBe}$   
 $\rightarrow \text{ab}bcbcb\text{Be}$   
 $\rightarrow \text{ab}bcbcbde$



## Classificação de *parsers*

Top-down parsers  
operam desta forma!

Número de tokens de  
lookahead lidos porém  
não consumidos

$LL(k)$

Lê entrada da esquerda  
(L) para a direita

Procura derivação mais à  
esquerda



## Resultados gerais

Pode-se construir parser LL( $k=1$ ) para uma gramática não-ambígua.

Gramáticas ambíguas pedem parsers com valor de  $k > 1$ .

Existem algoritmos para gerar parsers top-down de gramáticas ambíguas (e também com recursão à esquerda).



## Predictive parsing

- É um parser **top-down**
- É um parser que não requer *backtracking*
  - Simples de construir manualmente
  - Mas, requer modificação na gramática para tratar recursão à esquerda e ambigüidade
- Eliminação de recursão à esquerda (fatoração):

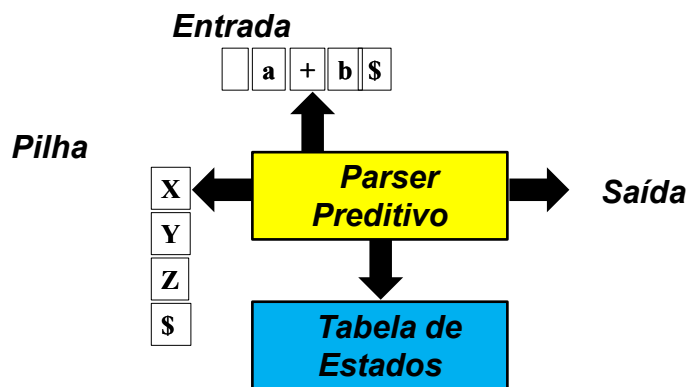
$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$



$S \rightarrow aABe$   
 $A \rightarrow bK$   
 $K \rightarrow bck \mid \epsilon$   
 $B \rightarrow d$



# Non-Recursive Predictive Parsing



## Bottom-up parsing

- Usa as produções de forma invertida até chegar no símbolo inicial.
- O *string* **abbcbcde** é caracterizado pela sequência de derivações abaixo.

$S \rightarrow aABe$   
 $A \rightarrow Abc \mid b$   
 $B \rightarrow d$

**a**bbcbcd**e**  
→ a**A**bbcbcd**e**  
→ a**A**bcde  
→ a**A**de  
→ a**A**Be  
→ S



## Shift-reduce parsers

- Usa uma pilha e uma tabela
  - Pilha: armazena tokens para salvar contexto
  - Tabela: determina as opções atuais de ação
- Possíveis ações:
  - Shift (push): coloca tokens na pilha
  - Reduce: determina uma produção
  - Accept: finaliza. reconhece string.
  - Error: Nenhuma ação é possível



## Shift-Reduce Parsers

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

Pilha	Entrada	Ação
\$	$id_1 + id_2 * id_3 \$$	shift
$\$id_1$	$+ id_2 * id_3 \$$	reduce by $E \rightarrow id$
$\$E$	$+ id_2 * id_3 \$$	shift
$\$E +$	$id_2 * id_3 \$$	shift
$\$E + id_2$	$* id_3 \$$	reduce by $E \rightarrow id$
$\$E + E$	$* id_3 \$$	shift
$\$E + E *$	$id_3 \$$	shift
$\$E + E * id_3$	$\$$	reduce by $E \rightarrow id$
$\$E + E * E$	$\$$	reduce by $E \rightarrow E * E$
$\$E + E$	$\$$	reduce by $E \rightarrow E + E$
$E$	$\$$	accept



## Conflitos Shift-Reduce

- ⌘ Há GLC onde parsers shift-reduce não podem ser usados.
- ⌘ Surgem conflitos e o parser não consegue decidir se shift ou reduce.
- ⌘ Tais gramáticas são ditas não LR.
- ⌘ Gramáticas ambíguas são não-LR.
- ⌘ A gramática com o else pendente não pode ser processada por parser LR.

## Parsers LR

- Eficientes
- Reconhecem quase todos os construtores de linguagens de programação expressos por GLC.
- O mais geral sem *backtracking*.
- A classe de gramáticas tratadas é um superconjunto próprio dos parsers preditivos.

## Parsers LR

- Pode detectar erros sintáticos o mais cedo possível.
- Difíceis de serem gerados manualmente!!!
- Facilmente automatizáveis a partir de GLC!
- Técnicas de construção:
  - SLR (Simple LR) - +Fácil implementar, - Poder
  - LR Canonico - +Custo implementação, +Poder
  - LALR (Look-ahead) – intermediário Custo, Poder



## Classificação de *parsers*

Bottom-up parsers  
operam desta forma!

Número de tokens de  
*lookahead* lidos porém não  
consumidos

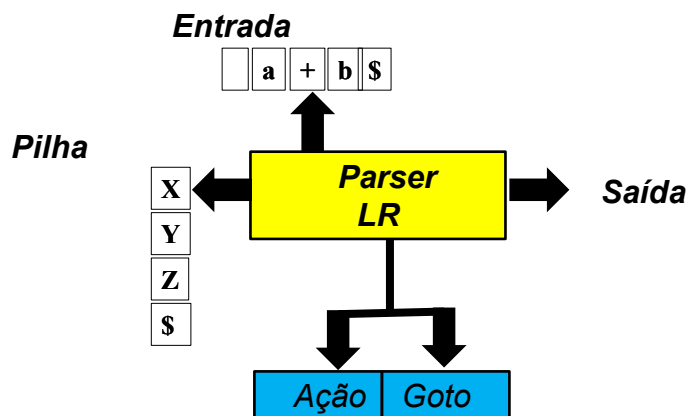
LR(k)

Lê entrada da esquerda  
(L) para direita

Procura a derivação  
mais à direita



# Non-Recursive Predictive Parsing



## Usando Gramáticas Ambíguas

### Precedência e associatividade dos operadores

A gramática ambígua:

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

Pode ser re-escrita como a não-ambígua:

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

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

$$F \rightarrow (E) \mid id$$

Onde + tem menor precedência que \*





# Usando Gramáticas Ambíguas

## Ambiguidade do “Dangling-else”

A gramática ambígua:

```
stmt → if expr then stmt  
      | if expr then stmt else stmt  
      | other
```

Deve-se sempre dar prioridade ao *shift* juntando o **else** ao primeiro **if\_then**.

# Usando Gramáticas Ambíguas

## Produções Especiais

A gramática ambígua:

$$E \rightarrow E^{\wedge} E\_E \mid E\_E \mid E^{\wedge} E \mid \{E\} \mid c$$
$$A^{\wedge} B\_C \rightarrow \boxed{A^B_C \quad A^{B_c}}$$

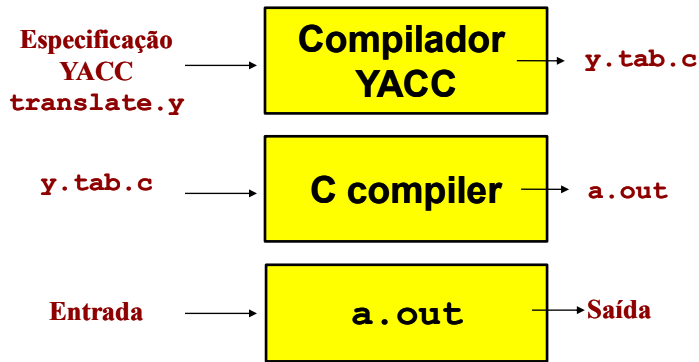
Não pode ser feita não-ambígua com precedência e associatividade.

Uma nova produção deve ser incorporada:

$$A^{\wedge} B\_C \rightarrow A^B_C$$

# YACC

## Yet Another Compiler Compiler



# YACC

## Yet Another Compiler Compiler

- Um programa YACC tem três partes:
  - Declarações  
%%
  - Regras de tradução  
%%
  - Rotinas auxiliares
- A Gramática de expressões aritméticas:
  - $E \rightarrow E + T \mid T$
  - $T \rightarrow T * F \mid F$
  - $F \rightarrow (E) \mid \text{digit}$



## Calculadora em YACC

$E \rightarrow E+T \mid T$

$T \rightarrow T * F \mid F$

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

```
%{
#include <ctype.h>
%}

%token DIGIT

%%
line : expr '\n' { printf("%d\n", $1); }
;
expr : expr '+' term { $$ = $1 + $3; }
    | term
;
term : term '*' factor { $$ = $1 * $3; }
    | factor
;
factor : '(' expr ')' { $$ = $2; }
      | DIGIT
;

%%
yylex() {
    int c;
    c = getchar();
    if (isdigit(c)) {
        yylval = c - '0';
        return DIGIT;
    }
    return c;
}
```

O pré-processador C inclui o predicado *isdigit*

### Declarações

### Regras de tradução

$E \rightarrow E+T \mid T$

$T \rightarrow T * F \mid F$

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

### Rotina auxiliar Chamada ao Lex

```
%{
#include <ctype.h>
%}

%token DIGIT

%%
line : expr '\n' { printf("%d\n", $1); }
;
expr : expr '+' term { $$ = $1 + $3; }
    | term
;
term : term '*' factor { $$ = $1 * $3; }
    | factor
;
factor : '(' expr ')' { $$ = $2; }
      | DIGIT
;

%%
yylex() {
    int c;
    c = getchar();
    if (isdigit(c)) {
        yylval = c - '0';
        return DIGIT;
    }
    return c;
}
```

O pré-processador C inclui o predicado *isdigit*

Ações semânticas

## Conflitos no YACC

Yet Another Compiler Compiler

- Conflitos *reduce/reduce* são resolvidos pela ordem em que aparecem as produções.
- Conflitos *shift/reduce* sempre resolvidos como *shift* (“*dangling-else*”).
- %left ‘+’ ‘-’ (mesma precedência e associativo à esquerda)
- %right ‘^’ (associativo à direita)
- %nonassoc ‘<’ (não-associativo)
- Tokens tem precedência crescente na ordem em que aparecem.



## Comparação Top-down e Bottom-up

- Em geral, bottom-up é mais poderoso
  - Coloca menos restrições na gramática
- Bem mais trabalhoso de se escrever e manter manualmente. Porém...
  - Yacc e Bison geram parser de gramática LALR(1), um subconjunto LR(1)