
Universidade Católica de Pelotas

Análise Sintática

André Rauber Du Bois
dubois@ucpel.tche.br

ANÁLISE SINTÁTICA

- Cada linguagem possui regras que descrevem a estrutura sintática dos programas bem formados
- A sintaxe das linguagens de programação pode ser descrita através de gramáticas livres de contexto ou pela notação BNF (Backus-Naur)
- Gramáticas fornecem uma especificação sintática precisa e fácil de entender

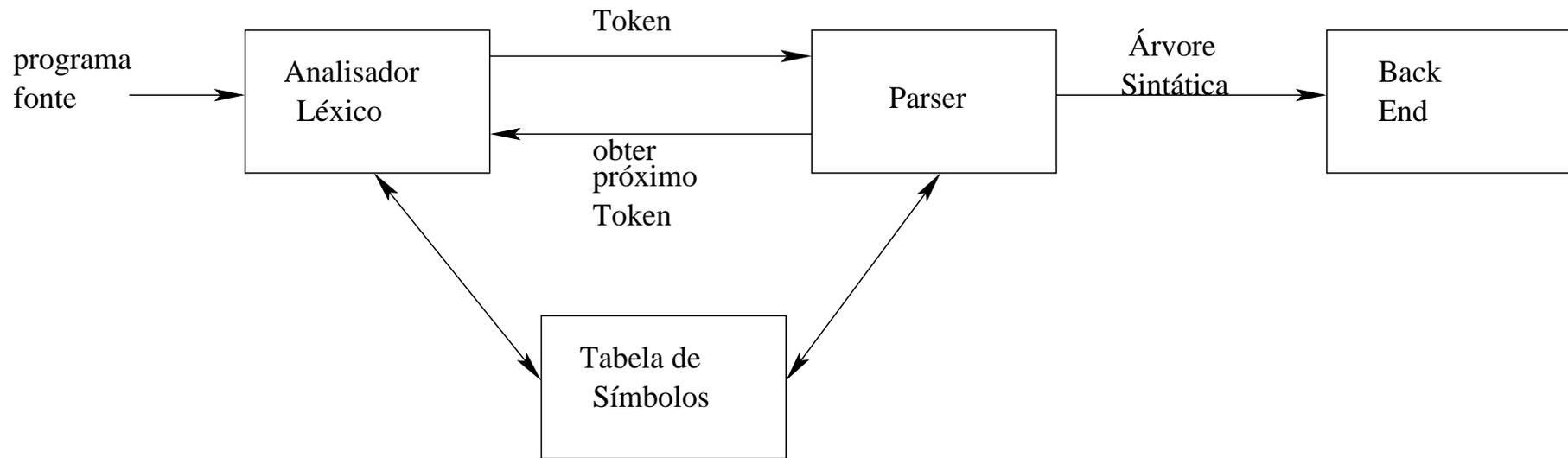
ANÁLISE SINTÁTICA

- Para certas gramáticas, pode-se construir automaticamente um Analisador Sintático que determine se o programa está sintaticamente bem formado
- O processo de construção do analisador pode detectar ambiguidades sintáticas
- Existem programas que geram parser partindo de gramáticas.
Ex: Yacc, javacc, etc.

ANÁLISE SINTÁTICA

- O analisador sintático (*parser*) obtém uma cadeia de *tokens* do analisador léxico e verifica se o programa é reconhecido pela gramática da linguagem fonte
- Os métodos mais usados para a implementação de analisadores sintáticos são o *bottom-up* e o *top-down*
- O *bottom-up*, constrói a árvore sintática de baixo para cima, ou seja, das folhas para a raiz
- O *top-down* contrói a árvore de cima para baixo, ou seja, da raiz para as folhas

O PARSER



ERROS DE SINTAXE

- Um bom compilador deve assistir o programador na identificação e correção de erros
- Os programas podem conter erros em muitos níveis
 - Léxicos: errar a grafia de um identificador ou palavra chave
 - Sintáticos: expressões aritméticas com parênteses não balanceados
 - Semânticos: operador aplicado a um operando incompatível
 - Lógicos: chamada infinitamente recursiva

ERROS DE SINTAXE

- Muitos dos erros são detectados durante a Análise sintática
- Erros geralmente ocorrem quando o programa desobedece a gramática da linguagem, ou seja, quando o fluxo de tokens vindo da análise léxica não respeita a gramática
- O tratador de erros do analisador sintático possui metas simples:
 - Deve relatar erros claramente
 - Deve se recuperar dos erros e continuar a análise para detectar erros subsequentes
 - Não deve retardar significativamente o processamento de programas corretos

ERROS DE SINTAXE

- Um erro é detectado quando a construção da árvore sintática falha
- Surgem então duas perguntas: que mensagem de erro deve ser fornecida ao usuário, e como prosseguir depois de um erro
- A posição em que o erro foi detectado pode não ter conexão com a posição real do erro que o usuário cometeu. Ex:

$$x = (p + q(- func(r - s) ;$$

- Provavelmente o erro do programa é que o parêntese depois do q deveria ser de fechamento, mas a maioria dos analisadores irá reportar a falta de um parêntese no final
- É quase impossível de se detectar o erro já que o q poderia ser uma função, e o sinal $-$ um menos monádico

ERROS DE SINTAXE

- Um sistema de detecção de erros simples poderia simplesmente parar quando um erro é encontrado
- Porém é frustrante para o programador ter que compilar várias vezes o mesmo programa, cada vez esperando que seja a última :-)
- Por isso programadores gostam de ter uma idéia de quantos erros restam no programa
- Existem basicamente duas estratégias de recuperação de erros, uma delas chamada de **correção de erros** e a outra **recuperação de erros sem correção**

CORREÇÃO E RECUPERAÇÃO SEM CORREÇÃO

- No sistema de **correção de erros** o fluxo de entrada é corrigido de alguma maneira para que a análise possa continuar
 - Problema: Uma correção no meio do programa pode ser uma correção errada, estragando a sintaxe de todo o resto do programa
- No sistema de **recuperação de erros sem correção** assim que um erro é encontrado, ele descarta toda a parte inicial do programa e recomeça a análise partindo do resto do programa
 - Problema: recuperação de erros confiável porém difícil de implementar já que a gramática deve ser modificada em tempo de execução
- A maioria dos geradores de analisadores sintáticos tem um

mecanismo interno de detecção e recuperação de erros, assim o projetista tem pouco a fazer sobre o assunto. Porém é importante saber como o tratamento de erros funcionada para tornar o compilador mais amigável

UM ANALISADOR SINTÁTICO TOP-DOWN

- Em uma gramática

$$N \rightarrow R \mid R1 \mid R2 \dots RN$$

- Dado um símbolo de entrada t um analisador top-down deve decidir qual alternativa R deve ser tomada para se reconhecer t
- Uma solução simples para o problema é usar-se funções recursivas booleanas que testem as alternativas de N em sequência
- Essa abordagem resulta em um *analisador sintático descendente recursivo*

UM ANALISADOR SINTÁTICO DESCENDENTE RECURSIVO

- Gramática:

entrada -> expressao EOF

expressao -> termo expressao_restante

termo -> IDENTIFICADOR | expressao_entre_parenteses

expressao_entre_parenteses -> '(' expressao ')'

expressao_restante -> '+' expressao | epsilon

- Gramática simples na qual o operador de + é associativo à direita

- Exemplo de expressão:

IDENTIFICADOR + (IDENTIFICADOR +
IDENTIFICADOR)

UM ANALISADOR SINTÁTICO DESCENDENTE RECURSIVO

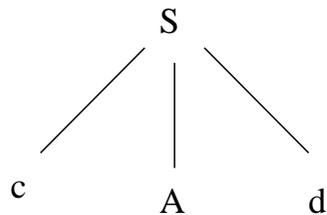
- Arquivo `parser.c`
- Cada regra N corresponde a uma rotina que retorna 1 (verdadeiro) se uma produção terminal de N foi encontrada. Quando encontrada a entrada é consumida
- Caso contrário a rotina retorna 0 (falso) e nenhuma entrada é consumida
- Os operadores lógicos (`&&`) e (`||`) são usados para requerer a presença de outros símbolos na entrada (`&&`), ou para representar uma alternativa (`||`)
- Se um membro de uma regra deve estar presente e não está, isso significa erro de sintaxe

PROBLEMAS DO AN. SINT. DESCENDENTE RECURSIVO

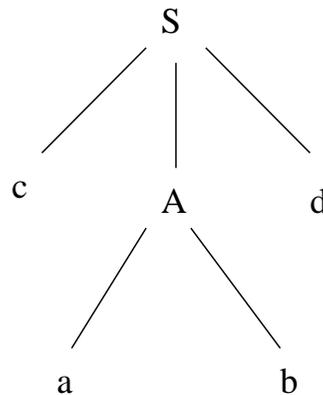
- O método geralmente exige retrocesso nas alternativas.
- Considere a seguinte gramática

$$S \rightarrow cAd$$
$$A \rightarrow ab \mid a$$

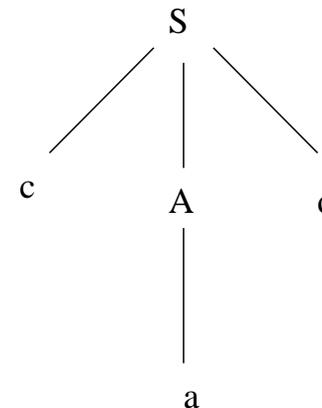
- e a cadeia de entrada $w = cad$



(1)



(2)



(3)

PROBLEMAS DO AN. SINT. DESCENDENTE RECURSIVO

- Primeiro, o apontador de entrada aponta para c , o primeiro caracter de entrada. A primeira regra S é escolhida e a árvore (1) é gerada
- O c é reconhecido pela folha à esquerda, e então a regra A é expandida gerando (2)
- O a é então reconhecido pela folha mais à esquerda, porém o próximo símbolo de entrada é d , e na árvore o próximo símbolo é o b
- Ocorre então o retrocesso, e a próxima alternativa de A é expandida gerando (3) e a entrada é reconhecida

PROBLEMAS DO AN. SINT. DESCENDENTE RECURSIVO

- Muitas vezes o método deixa de produzir um analisador correto
- Exemplo:

termo -> IDENTIFICADOR | elemento_indexado | ...

elemento_indexado -> IDENTIFICADOR '[' expressao ']'

- A sequência elemento_indexado nunca será experimentada, já que IDENTIFICADOR será consumido deixando o resto '[' expressao ']' de fora

PROBLEMAS DO AN. SINT. DESCENDENTE RECURSIVO

- Exemplo 2: Gramática para reconhecer ab e aab

$S \rightarrow A 'a' 'b'$

$A \rightarrow 'a' \mid \text{epsilon}$

```
int S(void) {  
    return (A() && require (token ('a')) &&  
           require (token('b')));  
}  
int A(void){  
    return token ('a') || 1;  
}
```

- ab não será reconhecido: A() consome 'a' e o require(token('a')) falha.

PROBLEMAS DO AN. SINT. DESCENDENTE RECURSIVO

- Exemplo 3: Recursão à esquerda

`expressao -> expressao '-' termo | termo`

```
int expressao(void) {  
    return (expressao () && require(token('-'))  
           && require (termo()) || require (termo()));  
}
```

- Loop infinito em `expressao()`

ELIMINAÇÃO DA RECURSÃO À ESQUERDA

- Uma gramática é recursiva à esquerda se existe um não terminal A e exista, direta ou indiretamente uma derivação $A \rightarrow Aa$
- Análises *top-down* não conseguem lidar com recursão à esquerda, por isso é necessário um método que a elimine
- Um par de produções recursivas á esquerda do tipo

$$A \rightarrow Aa \mid B$$

poderia ser substituído pelas produções não recursivas

$$A \rightarrow BA'$$

$$A' \rightarrow aA' \mid \text{epsilon}$$

ELIMINAÇÃO DA RECURSÃO À ESQUERDA

- Exemplo: Considere a seguinte gramática para expressões aritméticas

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

- Eliminado a recursão à esquerda, temos:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \text{epsilon} \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \text{epsilon} \\ F &\rightarrow (E) \mid \text{identificcador} \end{aligned}$$

ELIMINAÇÃO DA RECURSÃO À ESQUERDA

- Exemplo: Expressões aritméticas apresentada anteriormente

`expressao -> expressao '-' termo | termo`

- Eliminado a recursão à esquerda, temos:

`expressao -> termo expressao_final`

`expressao_final -> '-' termo expressao_final
| epsilon`

ELIMINAÇÃO DA RECURSÃO À ESQUERDA

- Cuidado: Recursão à esquerda indireta:

$$S \rightarrow AA \mid a$$
$$A \rightarrow SS \mid b$$

- Existem algoritmos para eliminar a recursão indireta

FATORAÇÃO À ESQUERDA

- É uma transformação que pode ser aplicada quando duas regras começam exatamente com o mesmo símbolo gramatical
- A idéia básica é postergar a decisão até quando tenhamos visto o suficiente da entrada para que possamos realizar a escolha certa
- Exemplo:

$$A \rightarrow aB1 \mid aB2$$

- Fatorando à esquerda obtemos:

$$A \rightarrow aA'$$

$$A' \rightarrow B1 \mid B2$$

FATORAÇÃO À ESQUERDA

- Exemplo:

```
cmd -> IF expr THEN cmd ELSE cmd  
      | IF expr THEN cmd
```

- Fatorando-se à esquerda:

```
cmd   -> IF expr THEN cmd cmd'  
cmd'  -> ELSE cmd | epsilon
```

FATORAÇÃO À ESQUERDA

- Exemplo do compilador Descendente Recursivo:

```
termo -> IDENTIFICADOR  
      | IDENTIFICADOR '[' expressao ']' | ...
```

- Fatorando-se à esquerda:

```
termo      -> IDENTIFICADOR id_depois | ...  
id_depois -> '[' expressao ']' | epsilon
```

SUBSTITUIÇÃO

- A substituição envolve substituir um não terminal N no lado direito de uma regra de produção pelas alternativas de N
- Ela é usada quando as entidades em conflito não são diretamente visíveis
- Exemplo

$$S \rightarrow p A q$$
$$A \rightarrow a \mid B c \mid \text{epsilon}$$

- Por substituição:

$$S \rightarrow p a q \mid p B c q \mid p q$$

SUBSTITUIÇÃO

- Exemplo

```
termo -> IDENTIFICADOR | elem_indexado | ..  
elem_indexado -> IDENTIFICADOR '[' expressao ']'
```

- Por substituição:

```
termo -> IDENTIFICADOR  
      | IDENTIFICADOR '[' expressao ']' | ...
```

- Agora temos um conflito direto que pode ser resolvido por fatoração à esquerda

SUBSTITUIÇÃO

- Problema apresentado anteriormente

$$S \rightarrow A \text{ 'a' 'b'}$$
$$A \rightarrow \text{'a'} \mid \text{epsilon}$$

- Por substituição:

$$S \rightarrow A \text{ 'a' 'b' } \mid \text{'a' 'b'}$$

ANÁLISE RECURSIVA PREDITIVA

- É possível implementar analisadores recursivos sem retrocesso
- A produção a ser aplicada é escolhida baseando-se no símbolo de entrada (análise preditiva)
- A gramática não pode ter recursividade à esquerda, deve estar fatorada à esquerda e para os não-terminais com mais de uma regra de produção, os primeiros terminais deriváveis devem ser capazes de identificar, sem equívocos, qual produção deve ser aplicada

ANÁLISE RECURSIVA PREDITIVA

- Dada uma entrada a e o não terminal A , deve-se saber qual das produções alternativas $A \rightarrow B1 \mid B2 \mid B3 \mid \dots$ deriva a sequência que inicia por a
- Temos então que calcular os conjuntos de primeiros símbolos produzidos por todas as alternativas na gramática, ou seja, os conjuntos FIRST
- O conjunto FIRST, de uma alternativa α , $FIRST(\alpha)$, contém todos os terminais com que α pode começar

ALGORITMO PARA CALCULAR O CONJUNTO FIRST

Inicialização:

1. se a é um terminal então $\text{FIRST}(a) = a$
2. Se $X \rightarrow \epsilon$ é uma produção, adicione ϵ em $\text{FIRST}(X)$

Regras de Inferência:

3. Para cada regra $N \rightarrow \alpha$, $\text{FIRST}(N)$ deve conter $\text{FIRST}(\alpha)$
4. Para cada alternativa α de forma $A\beta$, $\text{FIRST}(\alpha)$ deve conter $\text{FIRST}(A)$ menos o ϵ se este o contém
5. Para cada alternativa α de forma $A\beta$, e $\text{FIRST}(A)$ que contém ϵ , $\text{FIRST}(\alpha)$ deve conter $\text{FIRST}(\beta)$ inclusive ϵ se este o contém

EXEMPLO DE CONJUNTO FIRST

Dada a Gramática:

COMANDO \rightarrow CONDICIONAL | ITERATIVO | ATRIBUICAO

CONDICIONAL \rightarrow if EXPR then COMANDO

ITERATIVO \rightarrow repeat LISTA until EXPR |
while EXPR do COMANDO

ATRIBUICAO \rightarrow id := EXPR

- FIRST (CONDICIONAL) = if
- FIRST (ITERATIVO) = while, repeat
- FIRST (ATRIBUICAO) = id
- FIRST (COMANDO) = if, while, repeat, id

IMPLEMENTANDO A FUNCAO COMANDO

```
function COMANDO
  if token = 'if'
  then if CONDICIONAL
        then return true
        else return false
  else if token = 'while' or token = 'repeat'
  then if ITERATIVO
        then return true
        else return false
  else if token = 'id'
  then if ATRIBUICAO
        then return true
        else return false
  else return false
```

FUNÇÃO COMANDO

- Na função COMANDO agora não existe mais *backtracking*, ou *retrocesso*
- Cada alternativa é escolhida baseando-se no caracter de entrada
- Por isso este analisador é chamado de Analisador Recursivo Preditivo

PROBLEMAS COM O CONJUNTO FIRST

- Regras do tipo: $A \rightarrow \epsilon$, não podem ser reconhecidas
- Como a Regra A não inicia com nenhum símbolo, como podemos determinar se ela é a alternativa correta?
- Regra: Quando N produz um string não-vazio, vemos um token com o qual N pode se iniciar; quando N produz um string vazio, vemos o símbolo que pode *seguir* N
- Usa-se então o conjunto FOLLOW das alternativas

ALGORITMO PARA CALCULAR O CONJUNTO FOLLOW

1. Se s é o símbolo inicial e $\$$ representa fim de arquivo (EOF), então $\$$ está em FOLLOW(s)
2. Se existe produção do tipo $A \rightarrow \alpha X \beta$, então FIRST(β) está em FOLLOW(X)
3. Se existe produção do tipo $A \rightarrow \alpha X$ ou $A \rightarrow \alpha X \beta$, sendo que β gera ϵ , então todos os terminais que estiverem em FOLLOW(A) fazem parte de FOLLOW(X)

EXEMPLO: CALCULANDO FIRST

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \text{epsilon}$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \text{epsilon}$

$F \rightarrow -F \mid \text{id}$

- Começar sempre pelos não terminais mais fáceis. $\text{FIRST}(F) = \{-, \text{id}\}$
- $\text{FIRST}(T') = \{*, \epsilon\}$ e similarmente $\text{FIRST}(E') = \{+, \epsilon\}$
- Como T deriva FT' e F não deriva a palavra vazia, $\text{FIRST}(T) = \text{FIRST}(F) = \{-, \text{id}\}$
- Como E deriva TE' e T não deriva ϵ , $\text{FIRST}(E) = \text{FIRST}(T) = \{-, \text{id}\}$

EXEMPLO: CALCULANDO FOLLOW

- $\text{FOLLOW}(E) = \{\$\}$ pela regra 1, e $\text{FOLLOW}(E') = \text{FOLLOW}(E) = \{\$\}$, pois é o último símbolo do lado direito da regra $E \rightarrow TE'$
- partindo da produção $E' \rightarrow +TE'$, $\text{FOLLOW}(T)$ é obtido pela união de $\text{FIRST}(E')$ pela regra 2, com $\text{FOLLOW}(E')$ pela regra 3, já que $E' \rightarrow \epsilon$. Então $\text{FOLLOW}(T) = \{+,\$\}$
- $\text{FOLLOW}(T') = \text{FOLLOW}(T) = \{+,\$\}$, pois T' é o último na produção $T \rightarrow FT'$
- $\text{FOLLOW}(F)$ é obtido pela união entre $\text{FIRST}(T')$, pela regra 2, com $\text{FOLLOW}(T')$, pela regra 3, então $\text{FOLLOW}(F) = \{+,*,\$\}$

ANÁLISE PREDITIVA TABULAR

- É possível se construir analisadores preditivos *não recursivos* que utilizam uma pilha explícita, ao invés de chamadas recursivas
- Esse analisador implementa um autômato de pilha controlado por uma tabela de análise
- O analisador busca a produção a ser aplicada em uma tabela de análise, levando em conta o não-terminal no topo da pilha e o token sobre o cabeçote de leitura

ANÁLISE PREDITIVA TABULAR

- Exemplo, Gramática

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \text{epsilon}$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \text{epsilon}$

$F \rightarrow - F \mid \text{id}$

- E a tabela gerada usando os conjuntos FIRST and FOLLOW:

TABELA DE ANÁLISE

	id	+	*	-	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow -F$	

ANÁLISE PREDITIVA TABULAR

PILHA	ENTRADA	AÇÃO
\$E	id+id*id\$	E->TE'
\$E'T	id+id*id\$	T->FT'
\$ET'F	id+id*id\$	F->id
\$ET'id	id+id*id\$	desempilha e lê símbolo
\$E'T'	+id*id\$	T'->ε
\$E'	+id*id\$	E'->+TE'
\$E'T+	+id*id\$	desempilha e lê símbolo
\$E'T	id*id\$	T->FT'
\$E'T'F	id*id\$	F->id
\$E'T'id	id*id\$	desempilha e lê símbolo
\$E'T'	*id	T'->*FT'

ANÁLISE PREDITIVA TABULAR (CONTINUAÇÃO...)

PILHA	ENTRADA	AÇÃO
$\$E'T'F^*$	$*id\$$	desempilha e lê símbolo
$\$E'T'F$	$id\$$	$F \rightarrow id$
$\$E'T'id$	$id\$$	desempilha e lê símbolo
$\$E'T'$	$\$$	$T' \rightarrow \epsilon$
$\$E'$	$\$$	$E' \rightarrow \epsilon$
$\$$	$\$$	Aceita setença

ANÁLISE SINTÁTICA BOTTOM-UP

- A análise Bottom-up tenta construir uma árvore gramatical para uma cadeia de entrada começando pelas folhas, e trabalhando árvore acima em direção à raiz
- A cada passo de *redução*, uma subcadeia particular que reconheça o lado direito de uma produção, é substituída pelo símbolo à esquerda daquela produção
- Esse tipo de análise geralmente é implementada usando o método de empilhar e reduzir
- Vários geradores automáticos de parser são implementados dessa maneira, incluindo o famoso Yacc

ANÁLISE SINTÁTICA BOTTOM-UP

- Considerando a gramática

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

- A sentença *abbcd*e pode ser reduzida a S pelos seguintes passos:

$abbcd\bar{e} \rightarrow aA\bar{b}cd\bar{e} \rightarrow aA\bar{d}e \rightarrow aABe \rightarrow S$

IMPLEMENTAÇÃO DE PILHA DA AN. EMPILHAR E REDUZIR

- Uma forma conveniente de implementar um analisador sintático de empilhar e reduzir é usar uma pilha para guardar os símbolos gramaticais e um buffer de entrada para a cadeia w a ser decomposta

Pilha	Entrada
\$	w\$

- O analisador sintático empilha zero ou mais símbolos na pilha até que um *handle* (lado direito de uma produção) surja no topo da pilha. Este então é substituído por seu lado esquerdo. Este ciclo é repetido até termos uma entrada vazia e o símbolo inicial no topo da pilha

Pilha	Entrada
\$S	\$

ANALISADOR SINTÁTICO DE EMPILHAR E REDUZIR

- Considerando a seguinte gramática

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

ANALISADOR SINTÁTICO DE EMPILHAR E REDUZIR

PILHA	ENTRADA	AÇÃO
\$	id+id*id\$	empilhar
\$id	+id*id\$	reduzir
\$E	+id*id\$	empilhar
\$E+	id*id\$	empilhar
\$E+id	*id\$	reduzir
\$E+E	*id\$	reduzir
\$E	*id\$	empilhar
\$E*	id\$	empilhar
\$E*id	\$	reduzir
\$E*E	\$	reduzir
\$E	\$	aceitar