

Compiladores - Análise Recursiva

Fabio Mascarenhas – 2015.2

<http://www.dcc.ufrj.br/~fabiom/comp>

Geradores x Reconhecedores

- A definição formal de gramática dá um *gerador* para uma linguagem
- Para análise sintática, precisamos de um *reconhecedor*
- Mas podemos reformular a definição de gramática para dar um reconhecedor, também
- Uma PE-CFG (gramática livre de contexto com expressões de parsing) tem os mesmos conjuntos V , T e P de uma gramática tradicional, mas o conjunto P é uma função de não-terminais em expressões de parsing
- Podemos ter ou um não-terminal inicial S ou uma expressão de parsing inicial s

Expressões de Parsing

- Uma expressão de parsing é:
 - Um terminal a
 - Um não-terminal A
 - Uma *concatenação* de duas expressões pq
 - Uma *escolha* entre duas expressões p/q
 - A precedência da concatenação é maior que a da escolha, mas podemos usar parênteses para agrupamento
- 

Reconhecendo uma entrada

x, y, z são seqüências
de terminais
(ou ϵ)

- O significado de uma expressão de parsing p associada a uma gramática G , dada uma entrada qualquer, é dado por uma série de regras de dedução que dizem se a expressão reconhece um prefixo da entrada

$$\frac{}{G a ax \rightarrow x}$$

$$\frac{G P(A) xy \rightarrow y}{G A xy \rightarrow y}$$

$$\frac{G P xyz \rightarrow yz \quad G q yz \rightarrow z}{G Pq xyz \rightarrow z}$$

$$\frac{}{G \epsilon x \rightarrow x}$$

$$\frac{G P xy \rightarrow y}{G P/q xy \rightarrow y} \quad \text{esc. 1}$$

$$\frac{G q xy \rightarrow y}{G P/q xy \rightarrow y} \quad \text{esc. 2}$$

Exemplo

$$E \rightarrow TE'$$

$$E' \rightarrow (+TE' | -TE' | \epsilon)$$

$$T \rightarrow n | (\epsilon)$$

$$\int \begin{cases} E \rightarrow T \{ +T | -T \} \\ T \rightarrow n | (\epsilon) \end{cases}$$

$$\begin{array}{l} G \{ (n+n) \rightarrow n+n \} \dots \\ \hline G (E) (n+n) \rightarrow \epsilon \\ \hline G n | (E) (n+n) \rightarrow \epsilon \\ \hline G T (n+n) \rightarrow \epsilon \\ \hline G + + (n+n) \rightarrow (n+n) \end{array}$$

$$\underline{G n \quad n+(n+n) \rightarrow + (n+n)}$$

$$\underline{G n | (E) \quad n+(n+n) \rightarrow + (n+n)}$$

$$\underline{G T \quad n+(n+n) \rightarrow + (n+n)}$$

$$G + TE' + (n+n) \rightarrow \epsilon$$

$$\underline{G + TE' | -TE' | \epsilon \xrightarrow{+(n+n)} \epsilon}$$

$$\underline{G E' + (n+n) \rightarrow \epsilon}$$

$$\underline{G TE' \quad n+(n+n) \rightarrow \epsilon}$$

$$G E \quad n+(n+n) \rightarrow \epsilon$$

Repetição e Opcional

$$[P] \Rightarrow P/\epsilon \quad \{P\} \rightarrow P^* = P^* \cup \epsilon$$

- Regras de dedução podem para repetição e opcional podem ser dadas diretamente:

$$\frac{}{G P? X \rightarrow X} \stackrel{\text{op.1}}{=} \frac{G P X Y \rightarrow Y}{G P? X Y \rightarrow Y}$$

$$\frac{}{G P^* X \rightarrow X} \stackrel{\text{rep.1}}{=} \frac{G P X Y Z \rightarrow Y Z \quad G P^* Y Z \rightarrow Z}{G P^* X Y Z \rightarrow Z}$$

Não-determinismo da escolha

- As regras de dedução para a escolha não dizem qual das alternativas escolher: a escolha em uma gramática livre de contexto é *não-determinística*
- Simular não-determinismo em uma implementação real não é difícil, mas não é muito eficiente, e gera problemas de *ambiguidade*
- Todas as técnicas de análise sintática que vamos ver são diferentes maneiras de domar esse não-determinismo
- A primeira técnica, que vamos ver a seguir, reinterpreta a escolha para ser *determinística e ordenada*

Escolha ordenada

* exten

$$* \frac{a \neq b}{G \ a \ b \ x \rightarrow \text{falha}}$$

$$* \frac{G \ P(P) \rightarrow \text{falha}}{G \ A \ x \rightarrow \text{falha}}$$

$$* \frac{G \ P \ x \rightarrow \text{falha}}{G \ P \ q \ x \rightarrow \text{falha}}$$

$$* \frac{G \ P \ x \rightarrow x \quad G \ q \ y \rightarrow \text{falha}}{G \ P \ q \ x \ y \rightarrow \text{falha}}$$

$$\frac{G \ P \ x \ y \rightarrow \text{falha} \quad G \ q \ x \ y \rightarrow y \quad \text{esc. 2}}{G \ P \ q \ x \ y \rightarrow y}$$

$$* \frac{G \ P \ x \rightarrow \text{falha} \quad G \ q \ x \rightarrow \text{falha}}{G \ P \ q \ x \rightarrow \text{falha}}$$

Escolha ordenada (2)

$$\frac{G \ p \ x \rightarrow \text{falha}_{up-1}}{G \ p! \ x \rightarrow x}$$

$$\frac{G \ p \ x \rightarrow \text{falha}_{up-1}}{G \ p* \ x \rightarrow x}$$

Analizador Recursivo

- Maneira mais simples de implementar um analisador sintático a partir de uma gramática, mas não funciona com muitas gramáticas
- A ideia é manter a lista de tokens em um vetor, e o token atual é um índice nesse vetor
- Um **terminal** testa o token atual, e avança para o próximo token se o tipo for compatível, ou falha se não for
- Uma **sequência** testa cada termo da sequência, falhando caso qualquer um deles falhe
- Uma **alternativa** guarda o índice atual e testa a primeira opção, caso falhe volta para o índice guardado e testa a segunda, assim por diante

Analizador Recursivo

- Um **opcional** guarda o índice atual, e testa o seu termo, caso ele falhe volta para o índice guardado e não faz nada
- Uma **repetição** repete os seguintes passos até o seu termo falhar: guarda o índice atual e testa o seu termo
- Um **não-terminal** vira um procedimento separado, e executa o procedimento correspondente
- Construir a árvore sintática é um pouco mais complicado, as alternativas, opcionais e repetições devem jogar fora nós da parte que falhou!

Retrocesso Local *→ Escolhe determinística ordenada*

- Podemos definir o processo de construção de um parser recursivo com retrocesso local como uma transformação de EBNF para código Java
- Os parâmetros para nossa transformação são o termo EBNF que queremos transformar e um termo Java que nos dá o objeto da árvore sintática
- Vamos chamar nossa transformação de `$parser`
- `$parser(termo, arvore)` dá o código para análise sintática do termo, guardando o resultado em um ou mais nós de arvore caso seja bem sucedido

Retrocesso Local

- `$parser(termo, arvore)` dá o código para análise sintática do termo, guardando o resultado em um ou mais nós de `arvore` caso seja bem sucedido

```
$parser(terminal, arvore) =  
    ($arvore).child(match($terminal));
```

```
$parser(t1...tn, arvore) =  
    $parser(t1, arvore)  
    ...  
    $parser(tn, arvore)
```

```
$parser(NAOTERM, arvore) =  
    ($arvore).child(NAOTERM());
```

Retrocesso Local

- `$parser(termo, arvore)` dá o código para análise sintática do termo, guardando o resultado em um ou mais nós de arvore caso seja bem sucedido

```
$parser(t1 | t2, arvore) =  
{  
    int atual = pos;  
    try {  
        Tree rascunho = new Tree();  
        $parser(t1, rascunho);  
        ($arvore).children.addAll(rascunho.children);  
    } catch (Falha f) {  
        pos = atual;  
        $parser(t2, arvore);  
    }  
}
```

$t_1 | (t_2 | t_3)$

Retrocesso Local

- `$parser(termo, arvore)` dá o código para análise sintática do termo, guardando o resultado em um ou mais nós de `arvore` caso seja bem sucedido

```
$parser([ termo ], arvore) =  
{  
    int atual = pos;  
    try {  
        Tree rascunho = new Tree();  
        $parser(termo, rascunho);  
        ($arvore).children.addAll(rascunho.children);  
    } catch(Falha f) {  
        pos = atual;  
    }  
}
```

Retrocesso Local

- `$parser(termo, arvore)` dá o código para análise sintática do termo, guardando o resultado em um ou mais nós de `arvore` caso seja bem sucedido

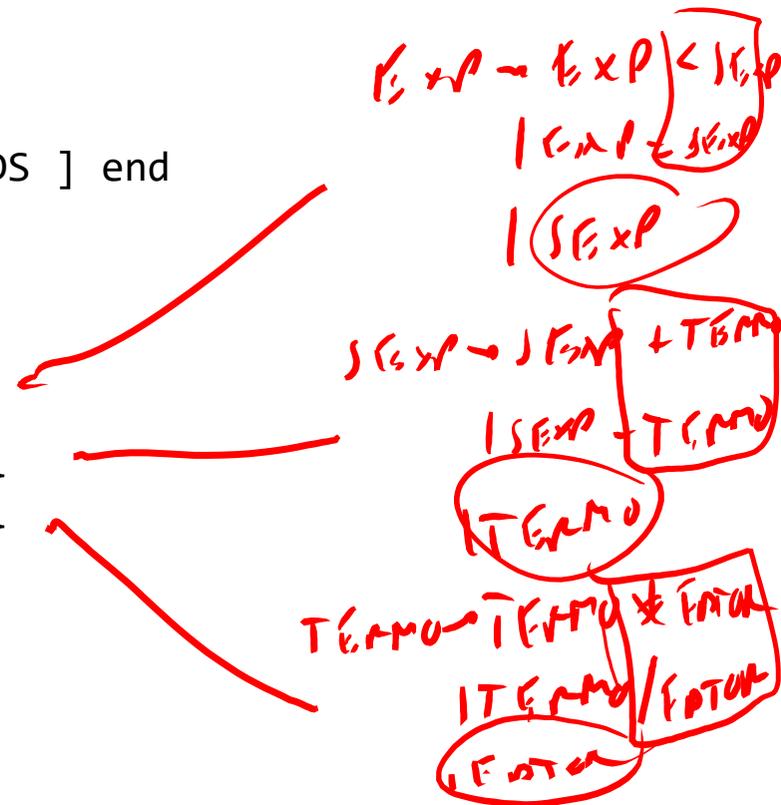
```
$parser({ termo }, arvore) =  
  while(true) {  
    int atual = pos;  
    try {  
      Tree rascunho = new Tree();  
      $parser(termo, rascunho);  
      ($arvore).children.addAll(rascunho.children);  
    } catch(Falha f) {  
      pos = atual;  
      break;  
    }  
  }  
}
```

Um analisador recursivo para TINY

com a trilha

- Vamos construir um analisador recursivo para TINY de maneira sistemática, gerando uma árvore sintática
- O vetor de tokens vai ser gerado a partir de um analisador léxico escrito com o JFlex

```
S      -> CMDS
CMDS   -> CMD { ; CMD }
CMD    -> if EXP then CMDS [ else CMDS ] end
        | repeat CMDS until EXP
        | id := EXP
        | read id
        | write EXP
EXP    -> SEXP { < SEXP | = SEXP }
SEXP   -> TERMO { + TERMO | - TERMO }
TERMO  -> FATOR { * FATOR | / FATOR }
FATOR  -> "(" EXP ")" | num | id
```



Detecção de erros

- Um analisador recursivo com retrocesso tem um comportamento ruim na presença de erros sintáticos
- Ele não consegue distinguir *falhas* (um sinal de que ele tem que tentar outra possibilidade) de *erros* (o programa está sintaticamente incorreto)
- Uma heurística é manter em uma variável global uma marca d'água que indica o quão longe fomos na sequência de tokens

Retrocesso local x global

- O retrocesso em caso de falha do nosso analisador é *local*. Isso quer dizer que se eu tiver $(A \mid B) C$ e A não falha mas depois C falha, ele não tenta B depois C novamente

- Da mesma forma, se eu tenho $A \mid AB$ a segunda alternativa nunca vai ser bem sucedida

$i \vdash (B \vee A) \text{ Then } (C \vee D)$

- As alternativas precisam ser *exclusivas*

$t_1 \mid t_2 \quad L(t_1) \cap L(t_2) \neq \emptyset$

$i \vdash (B \vee A) \text{ Then } (C \vee D) \text{ or } (C \vee D)$

- Retrocesso local também faz a repetição ser *gulosa*

\downarrow
 $\{ t_1 \} t_2$

- Uma implementação com retrocesso *global* é possível, mas mais complicada

Recursão à esquerda

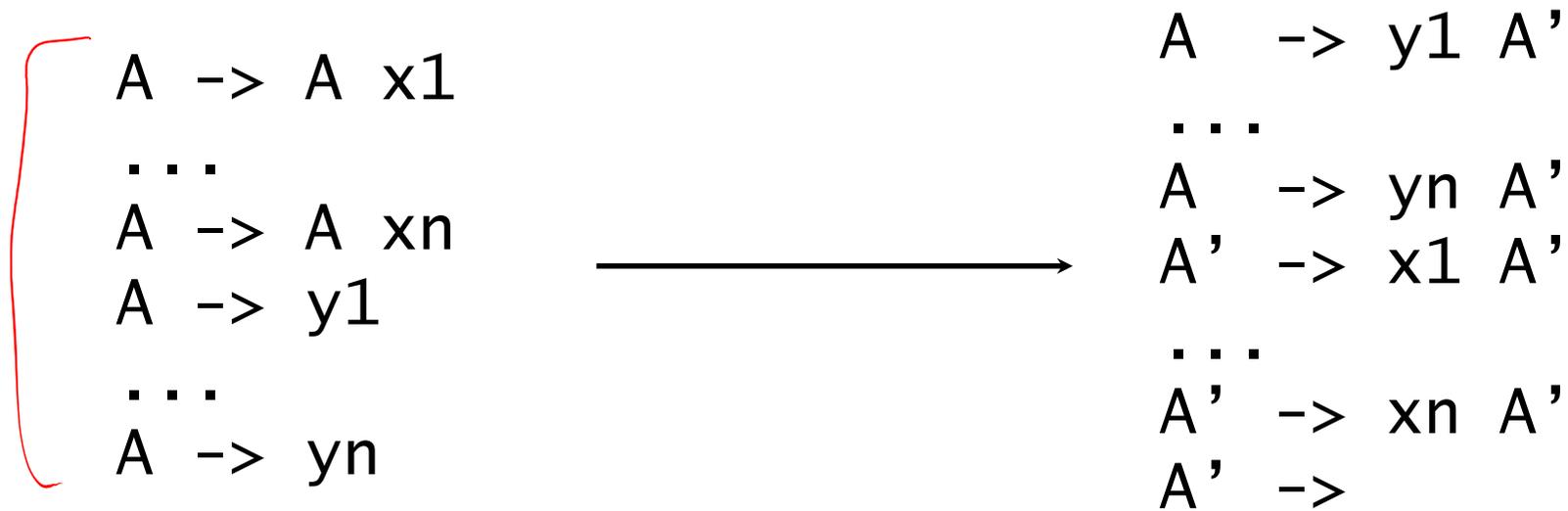
- Outra grande limitação dos analisadores recursivos é que as suas gramáticas não podem ter *recursão à esquerda*
- A presença de recursão à esquerda faz o analisador entrar em um laço infinito!
- Precisamos transformar recursão à esquerda em repetição
- Fácil quando a recursão é direta:

$$A \rightarrow A \text{ (x1) } | \dots | A \text{ (xn) } | y1 | \dots | yn$$

↓

$$A \rightarrow (y1 | \dots | yn) \{ x1 | \dots | xn \}$$

Eliminação de recursão sem EBNF



Parsing Expression Grammars

- As parsing expression grammars (PEGs) são uma generalização do parser com retrocesso local
- A sintaxe das gramáticas adota algumas características de expressões regulares: * e + para repetição ao invés de {}, ? para opcional ao invés de []
- Usa-se / para alternativas ao invés de |, para enfatizar que esse é um operador bem diferente do das gramáticas livres de contexto
- Acrescentam-se dois operadores de *lookahead*: &e e !e *La positivo*
- Finalmente, uma PEG pode misturar a tokenização com a análise sintática, então os terminais são *caracteres* (com sintaxe para strings e classes)

Uma PEG para TINY

if <<2

if <<2

```
S      <- CMDS
CMDS   <- CMD (; CMD)*
CMD    <- if EXP then CMDS (else CMDS)? end
        / repeat CMDS until EXP
        / id := EXP
        / read id
        / write EXP
EXP    <- SEXP (< SEXP / = SEXP)*
SEXP   <- TERMO (+ TERMO / - TERMO)*
TERMO  <- FATOR ("*" FATOR / "/" FATOR)*
FATOR  <- "(" EXP ")" / id / num
```

id < SP " " "+" ! [a-zA-Z0-9_]

1,0 +