

Compiladores

Análise Sintática

Introdução

Profa. Dra. Ana Carolina Lorena
UNIFESP

Análise Sintática

- Tarefa do analisador sintático é **determinar a estrutura sintática** de um programa a partir das marcas produzidas pela varredura (análise léxica)
 - ℓ Constrói no processo, explícita ou implicitamente, uma árvore de análise sintática, que representa essa estrutura

sequência de marcas  árvore sintática

analisador sintático

- ℓ As marcas são solicitadas pelo analisador sintático (Parser) pela ativação de `getToken`

Análise Sintática

- Um problema mais difícil é o **tratamento de erros**
 - ℓ O analisador sintático precisa registrar uma mensagem de erro e também se recuperar a condição sem erros e seguir em frente (para descobrir o maior número de erros possível)
 - ℓ Por vezes, ele pode efetuar **correção de erros** (normalmente só em casos simples)

Sintaxe

- **Sintaxe de uma linguagem de programação:**
normalmente especificada por regras gramaticais de uma **gramática livre de contexto**
 - ▮ Maior diferença em relação a gramáticas regulares de análise léxica é que regras sintáticas podem ser **recursivas**
 - ▮ Ex. if pode ser aninhado com outros ifs, o que exige recursão
 - ▮ Ex de regra: $\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{número}$
 - ▮ $\text{op} \rightarrow + \mid - \mid *$

Gramáticas Livres de Contexto

- **Definição:** $G = (T, N, P, S)$, em que
 - ▮ T é um **conjunto de terminais**
 - ℳ Terminam derivações
 - ▮ N é um **conjunto de não-terminais**, disjunto de T
 - ℳ Não terminam derivações, sempre podem ser substituídos
 - ▮ P é um **conjunto de produções** ou regras gramaticais, na forma $A \rightarrow \alpha$, em que $A \in N$ e $\alpha \in (T \cup N)^*$
 - ℳ Produzem cadeias
 - ▮ $S \in N$ é um **símbolo inicial**
- ▮ Conjunto $T \cup N$ é denominado **conjunto de símbolos** de G
- ▮ Cadeia $\alpha \in (T \cup N)^*$ é denominada **forma sentencial**

Gramáticas Livres de Contexto

Passo de derivação sobre G: tem a forma $\alpha A \gamma \Rightarrow \alpha \beta \gamma$, em que:

$$\alpha, \beta \text{ e } \gamma \in (T \cup N)^*$$

$$A \rightarrow \beta \in P$$

- **Derivação sobre G:** tem a forma $S \Rightarrow^* w$, em que:
 - ▮ $w \in T^*$ (cadeia apenas com terminais ou **sentença**)
 - ▮ Uma **derivação** é uma sequência de substituições de não-terminais por escolhas à direita das regras gramaticais
 - ▮ A cada passo da derivação, uma única substituição é feita com base em uma escolha de regra gramatical
 - ▮ Regras **definem** e os passos de derivação **constroem** por substituição

Gramáticas Livres de Contexto

- **Linguagem gerada por G :** $L(G) = \{w \in T^* \mid \exists \text{ uma derivação } S \Rightarrow^* w \text{ de } G\}$
É o conjunto de **sentenças deriváveis de G**
- **Derivação à esquerda:** $S \Rightarrow^*_{lm} w$ é uma derivação em que cada passo $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ é tal que $\alpha \in T^*$ (tem apenas terminais)
- **Derivação à direita:** $S \Rightarrow^*_{rm} w$ é uma derivação em que cada passo $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ é tal que $\gamma \in T^*$

Gramáticas Livres de Contexto

- Um conjunto de cadeias L é denominado **linguagem livre de contexto** se existir uma gramática livre de contexto G tal que $L = L(G)$

Exemplos

- ℓ) G com a regra $E \rightarrow (E) \mid a$
Tem um não-terminal E e três terminais (,) e a
Gera linguagem $L(G) = \{a, (a), ((a)), \dots\} = ({}^n a)^n$
- ℓ) G com regra $E \rightarrow (E)$ não gera cadeia alguma
 $L(G) = \{\}$
Não possui caso base para a recursão e fica em recursão infinita
- ℓ) G com regra $E \rightarrow E+a \mid a$
Exercício 1: a) quais são os terminais e não-terminais?
b) qual a linguagem gerada?

Gramáticas Livres de Contexto

- Pela hierarquia de Chomsky, as Gramáticas Livres de Contexto são mais gerais e incluem as Gramáticas Regulares
 - ℓ Daria então para projetar um analisador sintático para ler diretamente caracteres do arquivo fonte e eliminar o sistema de varredura, mas isso comprometeria eficiência

Gramáticas Livres de Contexto

- Nome livre de contexto: não-terminais aparecem livres à esquerda da seta
 - ℓ $A \rightarrow \alpha$ indica que A pode ser substituída por α em *qualquer ponto*, independentemente de onde ocorra A
- Há requisitos em linguagens de programação que poderiam ser mais facilmente especificados por **gramáticas sensíveis ao contexto**, que são mais poderosas
 - ℓ Possui regras $\beta A \gamma \rightarrow \beta \alpha \gamma$, em que $\alpha \neq \varepsilon$ (regra se aplica apenas se β ocorre antes e γ ocorre depois do não-terminal)

Gramáticas Livres de Contexto

- Exemplo em que gramática sensível a contexto se aplicaria: regra de declaração de variável antes de uso

```
{ int x;
```

```
...
```

```
... x ...
```

```
}
```

(nome deve aparecer em declaração antes de ser usado)

- Problema é gramáticas sensíveis a contexto são mais difíceis de serem usadas como base para um analisador sintático
 - ℓ Isso é então deixado para a **análise semântica**

Gramáticas Livres de Contexto

- Nos compiladores:
 - ℓ Os **símbolos terminais** serão as **marcas** produzidas pelo sistema de varredura
 - ℓ Os **não-terminais** são nomes de **estruturas sintáticas**

Especificação BNF

- Forma de Backus-Naur (BNF) é normalmente usada para especificar as regras gramaticais das linguagens de programação
 - ℓ Possui meta-símbolo de escolha (|)
 - ℓ Concatenação também é usada
 - ℓ Não possui meta-símbolo de repetição, que é obtida pela recursão
 - ℓ Uso de seta \rightarrow para expressar definições de nomes
 - ℓ Podem usar expressões regulares como componentes

Não há padrão universal para convenções, vamos adotar uma comumente utilizada

Especificação BNF

- Regra gramatical em BNF define a estrutura cujo nome está à esquerda da seta

Exemplo: $\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{número}$
 $\text{op} \rightarrow + \mid - \mid *$

Primeira regra (recursiva) define uma estrutura de expressão
Segunda regra define um operador

Especificação BNF

- Repetição pode ser construída por recursão
 - ℓ Exemplo: $A \rightarrow Aa \mid a$ ou $A \rightarrow aA \mid a$
 - ℓ geram a linguagem $\{a^n \mid n \geq 1\}$

Primeira regra ($A \rightarrow Aa \mid a$) é **recursiva à esquerda**, pois não-terminal A aparece como primeiro símbolo à direita da regra que define A

Segunda regra ($A \rightarrow aA \mid a$) é **recursiva à direita**

Especificação BNF

- Considere uma regra recursiva à esquerda $A \rightarrow A\alpha|\beta$
 - ℓ α e β representam cadeias arbitrárias e β não começa com A
 - ℓ Ela gera todas as cadeias da forma $\beta, \beta\alpha, \beta\alpha\alpha, \dots$
- Considere uma regra recursiva à direita $A \rightarrow \alpha A|\beta$
 - ℓ α e β representam cadeias arbitrárias e β não termina com A
 - ℓ Ela gera todas as cadeias da forma $\beta, \alpha\beta, \alpha\alpha\beta, \dots$
- ℓ β também pode ser a cadeia vazia ε

Derivações e linguagem

- As regras da gramática definem o conjunto de cadeias de símbolos para marcas consideradas sintaticamente legais
- Exemplo: $\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{número}$
 $\text{op} \rightarrow + \mid - \mid *$

Expressão aritmética $(34-3)*42$ corresponde à cadeia legal com sete marcas: (número – número) * número

Já a cadeia $34-3*42$ não é uma expressão legal

Derivações e linguagem

- Regras da gramática determinam cadeias legais de símbolos para marcas fazendo uso de derivações
 - ℓ Uma derivação começa com um único nome de estrutura e termina com uma cadeia de símbolos para marcas
 - ℓ Exemplo: $\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{número}$
 $\text{op} \rightarrow + \mid - \mid *$

Expressão $(34-3)*42$

- (1) $\text{exp} \Rightarrow \text{exp op exp}$
- (2) $\Rightarrow \text{exp op número}$
- (3) $\Rightarrow \text{exp} * \text{número}$
- (4) $\Rightarrow (\text{exp}) * \text{número}$
- (5) $\Rightarrow (\text{exp op exp}) * \text{número}$
- (6) $\Rightarrow (\text{exp op número}) * \text{número}$
- (7) $\Rightarrow (\text{exp} - \text{número}) * \text{número}$
- (8) $\Rightarrow (\text{número} - \text{número}) * \text{número}$

Derivações e linguagem

- Conjunto de cadeias de símbolos para marcas obtido por derivações para exp é a linguagem definida pela gramática de expressões
 - ℓ Contém todas as expressões sintaticamente legais
 - ℓ $L(G) = \{s \mid \text{exp} \Rightarrow^* s\}$, em que G representa a gramática de expressões

Exemplos

◦ Gramática de declarações simplificada

`declaração → if-decl | outra`

`if-decl → if (exp) declaração |
 if (exp) declaração else declaração`

`exp → 0 | 1`

Linguagem = declarações if aninhadas

Exemplos: `if (0) outra`

`if (1) outra else outra`

`if (1) outra else if (0) outra else outra`

Exemplos

Gramática de declarações simplificada (versão 2)

`declaração` \rightarrow `if-decl` | `outra`

`if-decl` \rightarrow `if` (`exp`) `declaração` `else-parte`

`else-parte` \rightarrow `else` `declaração` | ϵ

`exp` \rightarrow `0` | `1`

ϵ -produção indica que estrutura `else-parte` é opcional

Exercício 2: mostre como é a derivação das seguintes sentenças:

a) `if` (`0`) `outra`

b) `if` (`1`) `outra` `else` `outra`

c) `if` (`1`) `outra` `else` `if` (`0`) `outra` `else` `outra`

Exemplos

- Gramática de parênteses balanceados
 $A \rightarrow (A) A \mid \varepsilon$
- Gramática para uma sequência de declarações
 $\text{decl-sequencia} \rightarrow \text{decl} ; \text{decl-sequencia} \mid \text{decl}$
 $\text{decl} \rightarrow s$
- Gera sequências com uma ou mais declarações separadas por ;
- $L(G) = \{s, s;s, s;s;s, \dots\}$
- Exercício 3:** o que ocorre se permitirmos a cadeia vazia na primeira regra anterior?
 $\text{decl-sequencia} \rightarrow \text{decl} ; \text{decl-sequencia} \mid \varepsilon$

Árvore de Análise Sintática

- Derivações não representam unicamente a estrutura das cadeias construídas
- Existem em geral muitas derivações para a mesma cadeia

Exemplo: expressão $(34-3)*42$

- (1) $\text{exp} \Rightarrow \text{exp op exp}$
- (2) $\Rightarrow (\text{exp}) \text{ op exp}$
- (3) $\Rightarrow (\text{exp op exp}) \text{ op exp}$
- (4) $\Rightarrow (\text{número op exp}) \text{ op exp}$
- (5) $\Rightarrow (\text{número op número}) \text{ op exp}$
- (6) $\Rightarrow (\text{número} - \text{número}) \text{ op exp}$
- (7) $\Rightarrow (\text{número} - \text{número}) * \text{exp}$
- (8) $\Rightarrow (\text{número} - \text{número}) * \text{número}$

Árvore de Análise Sintática

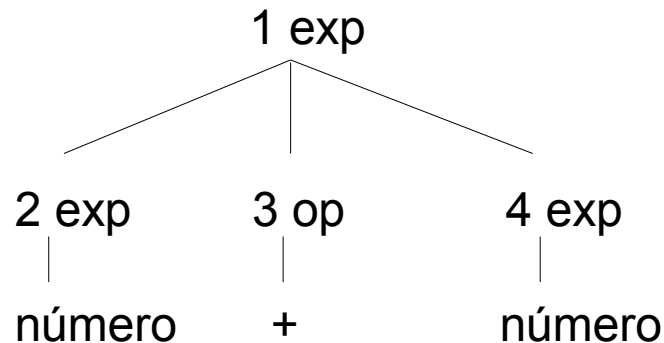
- ℓ Uma **árvore de análise sintática** correspondente a uma derivação é uma árvore com a raiz rotulada e com as propriedades:
 - ℓ Cada **nó** é rotulado com um **terminal**, um **não-terminal** ou ϵ
 - ℓ O **nó-raiz** é rotulado com o símbolo inicial S
 - ℓ Cada **nó folha** é rotulado com um **terminal** ou ϵ
 - ℓ Cada **nó não-folha** é rotulado com um **não-terminal**
 - ℓ Se um nó com rótulo $A \in N$ tiver n filhos com rótulos X_1, X_2, \dots, X_n (que podem ser terminais ou não-terminais), então $A \rightarrow X_1X_2\dots X_n \in P$

Árvore de Análise Sintática

- ℓ Cada derivação leva a uma árvore de análise sintática
- ℓ Em geral, muitas derivações podem levar à mesma árvore de análise sintática
- ℓ Cada árvore de análise sintática tem, entretanto, uma derivação à esquerda (percurso em pré-ordem) e uma à direita (inverso do percurso em pós-ordem) que são únicas e que levam a ela

Árvore de Análise Sintática

- Exemplo: (1) $\text{exp} \Rightarrow \text{exp op exp}$
(2) $\Rightarrow \text{número op exp}$
(3) $\Rightarrow \text{número} + \text{exp}$
(4) $\Rightarrow \text{número} + \text{número}$



Note que a enumeração está em pré-ordem

Árvore de Análise Sintática

- Exemplo: a mesma árvore corresponde à derivação:
exp \Rightarrow exp op exp
 \Rightarrow exp op número
 \Rightarrow exp + número
 \Rightarrow número + número
- Mas enumerações diferentes dos nós internos seriam aplicáveis

Árvore de Análise Sintática

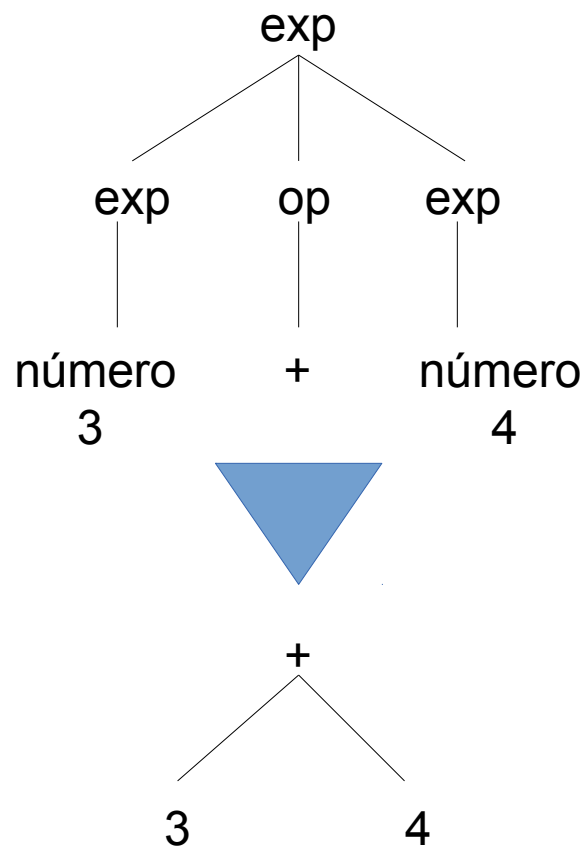
- ℓ A estrutura da árvore sintática depende da estrutura sintática específica da linguagem
 - ℓ Em geral é usada uma estrutura de dados dinâmica
 - ℓ Nós possuem campos para atributos

Árvore Sintática Abstrata

- ℓ Árvore de análise sintática é uma representação útil da estrutura de uma cadeia de marcas, que parecem como folhas
- ℓ Os nós internos representam passos de derivação
- ℓ Entretanto, ela contém mais informação do que o absolutamente necessário para um compilador
- ℓ **Árvores sintáticas abstratas** representam abstrações e contêm toda a informação necessária de forma mais eficiente
- ℓ São representação de uma notação simplificada, sintaxe abstrata

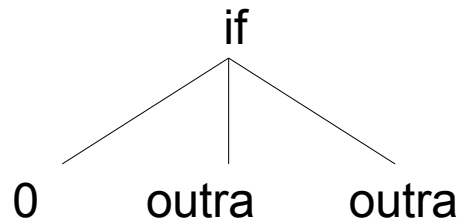
Árvore Sintática Abstrata

Exemplo 3 + 4

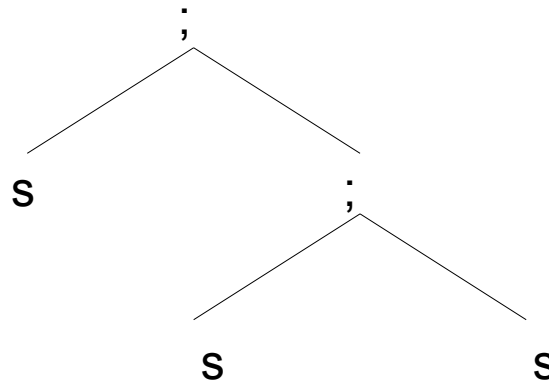


Árvore Sintática Abstrata

- Exemplo if (0) outra else outra



- Exemplo `s;s;s`



Árvore Sintática Abstrata

- Dependendo de como é construída a árvore sintática, podem ser definidas duas categorias de algoritmos de análise sintática:
 - Ascendentes
 - Descendentes

Análise Sintática

- Bibliografia consultada

Capítulo 3: LOUDEN, K. C. **Compiladores: princípios e práticas**. São Paulo: Pioneira Thompson Learning, 2004