

Construção de Compiladores

José de Oliveira Guimarães
Departamento de Computação
UFSCar - São Carlos, SP
Brasil
e-mail: jose@dc.ufscar.br

March 26, 2007

Contents

1	Introdução	2
1.1	Compiladores e Ligadores	2
1.2	Sistema de Tempo de Execução	4
1.3	Interpretadores	5
1.4	Aplicações	6
1.5	As Fases de Um Compilador	9
2	A Análise Sintática	12
2.1	Gramáticas	12
2.2	Ambigüidade	13
2.3	Associatividade e Precedência de Operadores	14
2.4	Modificando uma Gramática para Análise	15
2.5	Análise Sintática Descendente	17
2.6	Análise Sintática Descendente Recursiva	19
3	Análise Sintática Descendente Não Recursiva	23
3.1	Introdução	23
3.2	A Construção da Tabela M	25
3.3	Gramáticas LL(1)	29

Chapter 1

Introdução

1.1 Compiladores e Ligadores

Um compilador é um programa que lê um programa escrito em uma linguagem L_1 e o traduz para uma outra linguagem L_2 . Usualmente, L_1 é uma linguagem de alto nível como C++ ou Prolog e L_2 é assembler ou linguagem de máquina. Contudo, o uso de C como L_2 tem sido bastante utilizado. Utilizando C como linguagem destino torna o código compilado portátil para qualquer máquina que possua um compilador de C, que são praticamente 100% dos computadores. Se um compilador produzir código em assembler, o seu uso estará restrito ao computador que suporta aquela linguagem assembler.

Quando o compilador traduzir um arquivo (ou programa) em L_1 para linguagem de máquina, ele produzirá um arquivo de saída chamado de arquivo ou programa objeto, usualmente indicado pela extensão “.obj”. No caso geral, um programa executável é produzido pela combinação de vários arquivos objetos pelo ligador (*linker*). Veremos como o *linker* executa esta tarefa estudando um exemplo. Suponha que os arquivos “A.c” e “B.c” foram compilados para “A.obj” e “B.obj”. O arquivo “A.c” define uma função `f` e chama uma função `g` definida em “B.c”. O código de “B.c” define uma função `g` e chama a função `f`. Existe uma chamada a `f` em “A.c” e uma chamada a `g` em “B.c”. Esta configuração é ilustrada na Figura 1.1. O compilador compila “A.c” e “B.c” produzindo “A.obj” e “B.obj”, mostrados na Figura 1.2. Cada arquivo é representado por um retângulo dividido em três partes. A superior contém o código de máquina correspondente ao arquivo “A.c”. Neste código, utilizamos

```
call 000
```

para a chamada de *qualquer* função ao invés de colocarmos uma chamada de função em código de

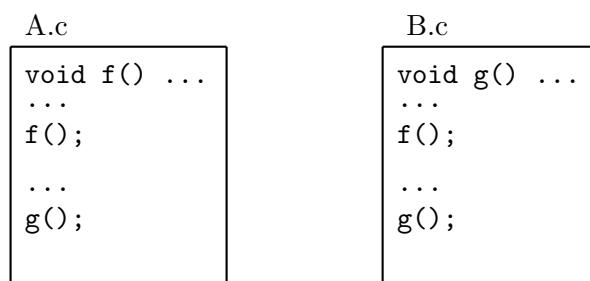


Figure 1.1: Dois arquivos em c que formam um programa

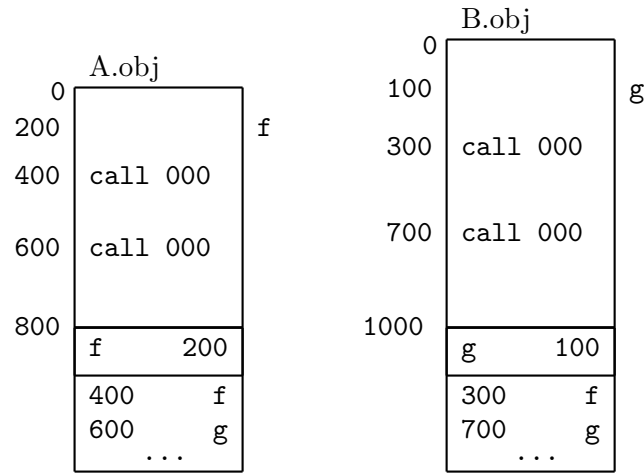


Figure 1.2: Configuração dos arquivos objeto

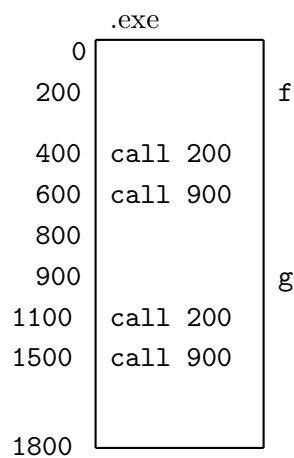


Figure 1.3: Configuração do arquivo executável

máquina, que seria composta apenas de números.

A parte intermediária do retângulo contém os nomes das funções definidas no arquivo juntamente com o endereço delas. Assim, o arquivo “A.obj” define uma função chamada *f* cujo endereço é 200. Isto é, o endereço da primeira instrução de *f* é 200. A última parte da definição de “A.obj”, o retângulo mais embaixo, contém os nomes das funções que são chamadas em “A.obj” juntamente com o endereço onde as funções são chamadas. Assim, a função *f* é chamada no endereço 400 e a função *g*, em 600. Os números utilizados acima (200, 400, 600) consideram que o primeiro byte de “A.obj” possui endereço 0.

Para construir o programa executável, o *linker* agrupa o código de máquina de “A.obj” e “B.obj” (parte superior do arquivo) em um único arquivo, mostrado na Figura 1.3. Como “B.obj” foi colocado após “A.obj”, o linker adiciona aos endereços de “B.obj” o tamanho de “A.obj”, que é 800. Assim, a definição da função *g* estava na posição 100 e agora está na posição 900 (100 + 800).

Como todas as funções estão em suas posições definitivas, o *linker* ajustou as chamadas das funções utilizando os endereços de *f* e *g*, que são 200 e 900. O arquivo “A.c” chama as funções *f* e *g*, como mostrado na Figura 1.1. O compilador transforma estas chamadas em código da forma

```

...
call 000 /* chama f */
...
call 000 /* chama g */
...

```

onde 000 foi empregado porque o compilador não sabe ainda o endereço de `f` e `g`. O *linker*, depois de calculado os endereços definitivos destas funções, modifica estas chamadas para

```

...
call 200 /* chama f */
...
call 900 /* chama g */
...

```

Ao executar o programa “.exe”, o sistema operacional carrega-o para a memória e pode ser necessário modificar todas as chamadas de função, adaptando-as para os endereços nos quais elas foram carregadas. Por exemplo, se o executável foi carregado para a posição 5000 na memória, a função que começa na posição 200 do executável estará na posição 5200 da memória. Assim, todas as chamadas a esta função, que são da forma

```
call 200
```

deverão ser modificadas para

```
call 5200
```

O carregamento do programa para a memória e possíveis realocações de endereços é chamada de carregamento.

1.2 Sistema de Tempo de Execução

Qualquer programa compilado requer um sistema de tempo de execução (*run-time system*). O sistema de tempo de execução (STE) é formado por todos os algoritmos de um programa executável que não foram especificados diretamente pelo programador. Algumas das funções do STE são descritas a seguir.

- Fazer a busca por um método em resposta a um envio de mensagem em C++.
- Procurar na pilha de funções ativas¹ qual função trata uma exceção levantada com `throw` em C++.
- Alocar memória para as variáveis locais de uma função antes do início da execução da mesma e desalocar a memória ao término da execução da função.
- Gerenciar o armazenamento do endereço de retorno da chamada da função.
- Chamar o construtor para uma variável cujo tipo é uma classe com construtor quando a variável for criada. Idem para o destrutor.
- Fazer a coleta de lixo de tempos em tempos. Todos os testes que o coletor de lixo possa ter inserido no programa fazem parte do sistema de tempo de execução.

¹As que foram chamadas pelo programa mas que ainda não terminaram a sua execução.

- Fornecer a *string* contendo a linha de comando com a qual o executável foi chamado.² A linha de comando fornecida pelo sistema operacional.
- Converter valores de um tipo para outro. Quando alguma transformação for necessária, será feita pelo STE.

Parte do código do sistema de tempo de execução é fornecido como bibliotecas já compiladas e parte é acrescentado pelo compilador. No primeiro caso, temos o algoritmo de coleta de lixo e o código que obtém a linha de comando do sistema operacional. No segundo caso estão todos os outros itens citados acima.

O *linker* agrupa todos os “.obj” que fazem parte do programa com um arquivo “.obj” que contém algumas rotinas do sistema de tempo de execução.

1.3 Interpretadores

Um compilador traduz o código fonte para linguagem de máquina que é modificada pelo *linker* para produzir um programa executável. Este programa é executável diretamente pelo computador.

Existe uma outra maneira de executar um programa: interpretá-lo. Há diversos modos de interpretação:

1. o interpretador lê o texto do programa e vai executando as instruções uma a uma. Atualmente esta forma é raríssima;
2. o interpretador toma o texto do programa e o traduz para uma estrutura de dados interna (semelhante à ASA que veremos neste curso, um conjunto de objetos que representa todo o programa). Então o interpretador, após a tradução do programa para a estrutura de dados, percorre esta estrutura interpretando o programa;
3. um compilador traduz o texto do programa para instruções de uma máquina virtual (pseudo-código). Então um interpretador executa estas instruções. Esta máquina virtual é usualmente uma máquina não só simples como feita sob medida para a linguagem que se quer utilizar

Existem vantagens e desvantagens no uso de compiladores/*linkers* em relação ao uso de interpretadores, sintetizados abaixo. Ao citar interpretadores, temos em mente aqueles descritos nos itens 2 e 3 acima.

1. O código interpretado pelo interpretador nunca interfere com outros programas. O interpretador pode conferir todos os acessos à memória e ao hardware impedindo operações ilegais.
2. Com o código interpretado torna-se fácil construir um *debugger* pois a interpretação está sob controle do interpretador e não do hardware.
3. Interpretadores são mais fáceis de fazer, por vários motivos. Primeiro, eles traduzem o código fonte para um pseudo-código mais simples do que o assembler ou linguagem de máquina utilizados pelos compiladores (em geral). Segundo, eles não necessitam de *linkers* — a ligação entre uma chamada da função e a função é feita dinamicamente, durante a interpretação. Terceiro, o sistema de tempo de execução do programa está presente no próprio interpretador, feito em uma linguagem de alto nível. O sistema de tempo de execução de um programa compilado deve ser, pelo menos em parte, codificado em assembler ou linguagem de máquina.

²Esta *string* pode ser manipulada em C++ com o uso dos parâmetros `argv` e `argc` da função `main`.

4. Interpretadores permitem iniciar a execução de um programa mais rapidamente do que se estivermos utilizando um compilador (veja a descrição 2 de interpretadores). Para compreender este ponto, considere um programa composto por vários arquivos já traduzidos para o pseudo-código pelo interpretador. Suponha que o programador faça uma pequena modificação em um dos arquivos e peça ao ambiente de programação para executar o programa. O interpretador traduzirá o arquivo modificado para pseudo-código e imediatamente iniciará a interpretação do programa. Apenas o arquivo modificado deve ser re-codificado. Isto contrasta com um ambiente de programação aonde compilação é utilizada.

Se um arquivo for modificado, freqüentemente será necessário compilar outros arquivos que dependem deste, mesmo se a modificação for minúscula. Por exemplo, se o programador modificar o valor de uma constante em um `define`

```
#define Num 10
```

de um arquivo “def.h”, o sistema deverá recompilar todos os arquivos que incluem “def.h”.

Depois de compilados um ou mais arquivos, dever-se-á fazer a ligação com o *linker*. Depois de obtido o programa executável, ele será carregado em memória e os seus endereços realocados para novas posições. Só após isto o programa poderá ser executado. O tempo total para realizar todas estas operações em um ambiente que emprega compilação é substancialmente maior do que se interpretação fosse utilizada.

5. Compiladores produzem código mais eficiente do que interpretadores. Como o código gerado por um compilador será executado pela própria máquina, ele será muito mais eficiente do que o pseudo-código interpretado equivalente. Usualmente, código compilado é 10-20 vezes mais rápido do que código interpretado.

Da comparação acima concluímos que interpretadores são melhores durante a fase de desenvolvimento de um programa, pois neste caso o importante é fazer o programa iniciar a sua execução o mais rápido possível após alguma alteração no código fonte.

Quando o programa estiver pronto, será importante que ele seja o mais rápido possível, o que pode ser obtido compilando-o.

1.4 Aplicações

Os algoritmos e técnicas empregados em compiladores são utilizados na construção de outras ferramentas descritas a seguir.

1. Editores de texto orientados pela sintaxe. Um editor de texto pode reconhecer a sintaxe de uma linguagem de programação e auxiliar o usuário durante a edição. Por exemplo, ele pode avisar que o usuário esqueceu de colocar (após o while:

```
while i > 0 )
```

2. *Pretty printers*. Um *pretty printer* lê um programa como entrada e produz como saída este mesmo programa com a tabulação correta. Por exemplo, se a entrada for

```
#include <iostream.h>
void
    main()
{ int i;      for ( i = 0;
```

```

        i < 10; i++
    )
    cout << endl;
}

```

a saída poderá ser

```

#include <iostream.h>

void main()
{
    int i;

    for ( i = 0; i < 10; i++ )
        cout << endl;
}

```

O comando `indent` do Unix formata e tabula um programa em C apropriadamente de acordo com algumas opções da linha de comando.

3. Analisadores estáticos de programas, que descubrem erros como variáveis não inicializadas, código que nunca será executado, situações em que uma rotina não retorna um valor, etc. O código abaixo mostra um erro que poderá ser descoberto por um destes analisadores.

```

void calcule( int *p, int *w )
{

    int soma = *p + *w;
    return soma + *p/*w; /* retorna soma */;
}

```

O Unix possui um analisador chamado `lint` que descobre erros deste tipo. Um analisador mais poderoso, disponível para o DOS e Windows, é o `PC-lint`.

4. Analisadores que retornam a cadeia de chamadas de um programa. Isto é, eles informam quem chama quem. Por exemplo, dado o programa

```

#include <stdio.h>

int fatorial( int n )
{
    if ( n >= 1 )
        return n *fatorial(n-1);
    else
        return 1;
}

int fat( int n )
{

```



```

    return n > 1 ? n*fatorial(n - 1) : 1 ;
}

void main()
{
    int n;

    printf("Digite n ");
    scanf( "%d", &n );
    printf("\nfat(%d) = %d\n", n, fat(n) );
}

```

como entrada ao utilitário cflow do Unix, a saída será:

```

1  main: void(), <lx.c 19>
2      printf: <>
3      scanf: <>
4      fat: int(), <lx.c 14>
5          fatorial: int(), <lx.c 6>
6          fatorial: 5

```

5. Formatadores de texto como $\text{T}_{\text{E}}\text{X}$ e $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$. Estes formatadores admitem textos em formato não documento com comandos especificando como a formatação deve ser feita e o tipo das letras. Por exemplo, o trecho

A função main inicia todos ... é utilizada com frequência para ...

é obtido digitando-se

A fun\c{c}\~{a}o {\tt main} inicia todos ... \'{e} utilizada com freq\{u}\^{}ncia para ...

6. Interpretadores de consulta a um Banco de Dados. O usuário pede uma consulta como

```
select dia > 5 and dia < 20 and idade > 25
```

que faz o interpretador imprimir no vídeo os registros satisfazendo às três condições acima.
7. Interpretadores de comandos de um sistema operacional. Além de reconhecer se o comando digitado é válido, o interpretador deve conferir os parâmetros:

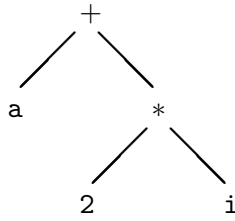
```

C:\>mkaedir so
C:\>dir *.cpp *.h
C:\>del *.exe

```

o interpretador sinalizaria erro nos dois primeiros comandos.

8. Interpretadores de expressões. Alguns programas gráficos leem uma função digitada pelo usuário e fazem o seu gráfico. A função é compilada para pseudo-código e que é então interpretado.

Figure 1.4: Árvore de sintaxe da sentença $a + 2*i$

9. Alguns editores de texto fazem busca por cadeias de caracteres utilizando expressões regulares. Por exemplo, para procurar por um identificador que começa com a letra A e termina com 0, podemos pedir ao editor para buscar por

```
A[a-zA-Z0-9_]*0
```

Um conjunto entre [e] significa qualquer caráter do conjunto e * é repetição de zero ou mais vezes do caráter anterior. Assim, [a-z]* é repetição de zero ou mais letras minúsculas. Esta expressão é transformada em um autômato e interpretada.

1.5 As Fases de Um Compilador

Um compilador típico divide a compilação em seis fases, descritas abaixo. Cada fase toma como entrada a saída produzida pela fase anterior. A primeira fase, análise léxica, toma como entrada o programa fonte a ser compilado. A última fase, geração de código, produz um arquivo “.obj” com o código gerado a partir do programa fonte.

1. Análise léxica. Esta fase lê o programa fonte produzindo como saída uma seqüência de números correspondentes a pequenas partes do programa chamados *tokens*. Por exemplo, às palavras chave `while`, `if` e `return` estão associados os números 1, 10, 12 e os números 40 e 45 estão associados a identificadores e literais numéricas. Assim, a entrada

```
if ( i > 5 )
    return j
```

produziria a seqüência de *tokens* 10, 50, 40, 60, 45, 51, 12, 40 e 66. Observe que 50, 60, 51 e 66 estão associados a “(”, “>”, “)” e “;”.

2. Análise sintática. Esta fase toma os números produzidos pela fase anterior e verifica se eles formam um programa correto de acordo com a gramática da linguagem. Em geral, uma árvore de sintaxe abstrata será também construída nesta fase. Esta árvore possuirá nós representando cada sentença da gramática. Por exemplo, a expressão $a + 2*i$ seria representada como mostrada na Figura 1.4.
3. Análise semântica. Esta fase confere se o programa está correto de acordo com a definição semântica da linguagem. Como exemplo, o comando

```
i = fat(n);
```

em C++ requer que:

- `i` e `n` tenham sido declarados como variáveis e `i` não tenha sido declarada com o qualificador `const`;
- `fat` seja uma função que possua um único parâmetro para cujo tipo o tipo de `n` possa ser convertido;

- `fat` retorne alguma coisa e que o tipo de retorno possa ser convertido para o tipo de `i`;

Todas estas conferências são feitas pelo analisador semântico. Observe que “`i = fat(n)`” é um comando sintaticamente correto, independente de qualquer conferência semântica.

4. Geração de código intermediário. Nesta fase, o compilador gera código para uma máquina abstrata que certamente é mais simples do que a máquina utilizada para a geração do código final. Este código intermediário é gerado para que algumas otimizações possam ser mais facilmente feitas. Aho, Sethi e Ullman citam o exemplo do comando

```
p = i + r*60
```

compilado para o código intermediário

```
temp1 = inttoreal(60)
temp2 = id3*temp1
temp3 = id2 + temp2
id1 = temp3
```

otimizado para

```
temp1 = id3*60.0
id1 = id2 + temp1
```

Note que:

- o tipo de `i` e `r` é real, sendo necessário converter `60` para real;
- `id1`, `id2` e `id3` correspondem a `p`, `i` e `r`;
- o código intermediário não admite, como assembler, gerar o código


```
id1 = id2 + id3*60
```

 diretamente. Tudo deve ser feito em partes pequenas;
- o compilador não gera o código otimizado diretamente. Primeiro ele gera um código ruim que então é otimizado. Isto é mais simples do que gerar o código otimizado diretamente.

5. Otimização de código. Como mostrado anteriormente, esta fase otimiza o código intermediário. Contudo, otimização de código pode ocorrer também durante a geração de código intermediário e na geração de código. Como exemplo do primeiro caso, o compilador pode otimizar

```
n = 20;
b = 4*k;
a = 3*n;
```

para

```
n = 20;
b = k << 2;
a = 60;
```

antes de gerar o código intermediário. Como exemplo do segundo caso, o compilador pode otimizar

```
i = i + 1
```

para

```
inc i
```

aproveitando o fato de que a máquina alvo possui uma instrução que incrementa o valor de um inteiro de um (`inc`). Admitimos que o código intermediário não possui esta instrução. Nesta fase também é feita a alocação dos registradores da máquina para as variáveis locais, o que também é um tipo de otimização.

6. Geração de código. Esta fase gera o código na linguagem alvo, geralmente assembler ou linguagem de máquina. Como visto na fase anterior, esta fase é também responsável por algumas otimizações de código.

As fases de um compilador são agrupadas em *front end* e *back end*. O *front end* é composto pelas fases que dependem apenas da linguagem fonte, como da sua sintaxe e semântica. O *back end* é composto pelas fases que dependem da máquina alvo.

Utilizando a divisão do compilador exposta anteriormente, o *front end* consiste da análise léxica, sintática e semântica, geração de código intermediário e otimização de código. O *back end* consiste apenas da geração de código. Obviamente, otimizações de código dependentes das características da máquina alvo fazem parte do *back end*.

Se tivermos um compilador para uma linguagem L que gera código para uma máquina A, poderemos fazer este compilador gerar código para uma máquina B modificando-se apenas o *back end*. Por este motivo, é importante produzir código intermediário para que nele possam ser feitas tantas otimizações quanto possíveis. Estas otimizações serão aproveitadas quando o compilador passar a gerar código para outras máquinas diferentes da original.

Em todas as fases, o compilador utilizará uma tabela de símbolos para obter informações sobre os identificadores do programa. Quando um identificador for encontrado durante a compilação, como `j` na declaração

```
int j;
```

ele será inserido na tabela de símbolos juntamente com outras informações, como o seu tipo. As informações da TS (tabela de símbolos) é utilizada para fazer conferências semânticas, como conferir que `j` foi declarado em

```
j = 1;
```

e que `1` pode ser convertido para o tipo de `j`. Ao gerar código para esta atribuição, o compilador consultará a TS para descobrir o tipo de `j` e então gerar o código apropriado.

Durante a análise léxica, todos os identificadores do programa estarão associados a um único número, que assumiremos ser 40. Assim,

```
j = i + 1;
```

retornará a seqüência de número 40 33 40 63 45. Para descobrir mais informações sobre o identificador, poderemos utilizar os métodos de um objeto `lex`. Por exemplo, para saber a *string* correspondente ao último identificador encontrado, fazemos

```
lex->getStrToken()
```

Com a *string* retornada, poderemos consultar a TS para obter mais informações sobre o identificador, como o seu escopo e tipo.

Nos próximos capítulos, veremos cada uma das fases do compilador em detalhes, embora com algumas diferenças das fases apresentadas nesta seção:

- o código gerado será em linguagem C;
- não haverá geração de código intermediário.

Chapter 2

A Análise Sintática

2.1 Gramáticas

A sintaxe de uma linguagem de programação descreve todos os programas válidos naquela linguagem e é descrita utilizando-se uma gramática, em geral livre de contexto. A gramática representa, de forma finita, todos os infinitos¹ programas válidos na linguagem.

Uma gramática G é especificada através de uma quádrupla (N, Σ, P, S) onde

- N é o conjunto de símbolos não-terminais;
- Σ é o conjunto de símbolos terminais;
- P é um conjunto de produções;
- S é o símbolo não-terminal inicial da gramática.

Os programas válidos da linguagem são aqueles obtidos pela expansão de S .

Utilizando a gramática

Expr	::= Expr "+" Term Expr "-" Term Term
Term	::= Term "*" Numero Term "/" Numero Numero
Numero	::= "0" "1" "2" "3" "4" "5" "6" "7" "8" "9"

temos que:

- $N = \{ \text{Expr, Term, Numero} \}$
- $\Sigma = \{ +, -, *, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$

Os terminais serão colocados entre aspas quando houver dúvidas sobre quais são os terminais. Neste exemplo, está claro que os operadores aritméticos e os dígitos são terminais e portanto as aspas não foram utilizadas.

- $P = \{ \text{Expr} ::= \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term}, \text{Term} ::= \text{Term} * \text{Numero} \mid \text{Term} / \text{Numero} \mid \text{Numero}, \text{Numero} ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \}$
- $S = \text{Expr}$

¹Em geral.

Nos exemplos deste curso, uma gramática será apresentada somente através de suas regras de produção com o símbolo inicial aparecendo do lado esquerdo da primeira regra. Assim, todos os elementos N , Σ , P e E estarão claramente identificados.

Uma sentença válida na linguagem é composta apenas por terminais e derivada a partir do símbolo inicial S , que neste exemplo é Expr . A derivação de Expr começa substituindo-se este símbolo pelos símbolos que aparecem à esquerda de Expr nas regras de produção da gramática. Por exemplo,

$$\text{Expr} \Rightarrow \text{Expr} + \text{Term}$$

é uma derivação de Expr . O lado direito de \Rightarrow pode ser derivado substituindo-se Expr ou Term . Substituiremos Expr :

$$\text{Expr} \Rightarrow \text{Expr} + \text{Term} \Rightarrow \text{Term} + \text{Term}$$

Substituindo-se sempre o não-terminal mais à esquerda, obteremos

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Expr} + \text{Term} \Rightarrow \text{Term} + \text{Term} \Rightarrow \text{Numero} + \text{Term} \\ &\Rightarrow 7 + \text{Term} \Rightarrow 7 + \text{Numero} \Rightarrow 7 + 2 \end{aligned}$$

Uma seqüência de símbolos w de tal forma que

$$S \Rightarrow \alpha_1 \Rightarrow \alpha_2 \dots \Rightarrow \alpha_n \Rightarrow w$$

$n \geq 0$, é chamada de sentença de G , onde G é a gramática (N, Σ, P, S) e w é composto apenas por terminais.

O símbolo $\xRightarrow{*}$ significa “derive em zero ou mais passos” e $\xRightarrow{+}$ significa “derive em um ou mais passos”. Então, uma sentença w de G é tal que

$$S \xRightarrow{+} w$$

A linguagem gerada pela gramática G é indicada por $L(G)$. Assim,

$$L(G) = \{w \mid S \xRightarrow{+} w\}$$

onde S é o símbolo inicial da gramática e w é composto apenas por terminais.

Uma derivação de S contendo símbolos terminais e não-terminais é chamada de *forma sentencial de G* . Isto é, se $S \xRightarrow{*} \alpha$, α é uma forma sentencial de G .

Neste curso, utilizaremos apenas derivações à esquerda e à direita. Em uma derivação à esquerda

$$S \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n \Rightarrow w$$

somente o não-terminal mais à esquerda de cada forma sentencial α_i ($1 \leq i \leq n$) é substituído a cada passo. Indicamos uma derivação à esquerda $\alpha \Rightarrow \beta$ por

$$\alpha \xrightarrow{lm} \beta$$

onde lm significa **leftmost**.

Derivação à direita possui uma definição análoga à derivação à esquerda, sendo indicada por

$$\alpha \xrightarrow{rm} \beta$$

onde rm significa **rightmost**.

Considerando a derivação

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Expr} + \text{Term} \Rightarrow \text{Term} + \text{Term} \Rightarrow \text{Numero} + \text{Term} \\ &\Rightarrow 3 + \text{Term} \Rightarrow 3 + \text{Numero} \Rightarrow 3 + 2 \end{aligned}$$

podemos construir uma árvore de análise associada a ela, mostrada na Figura 2.1. Esta árvore abstrai a ordem de substituição dos não-terminais na derivação de Expr : a mesma árvore poderia ter sido gerada por uma derivação à direita.

2.2 Ambigüidade

Uma gramática que produz mais de uma derivação à esquerda (ou direita) para uma sentença é dita ser ambígua. Isto é, uma gramática será ambígua se uma sentença w puder ser obtida a partir do símbolo inicial S por duas derivações à esquerda (ou direita) diferentes.

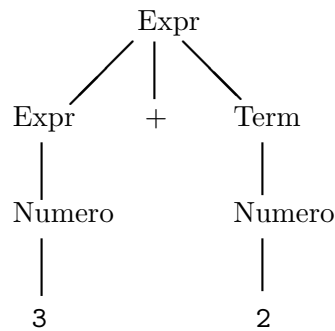


Figure 2.1: Árvore de análise de “3 + 2”

Uma gramática ambígua produz mais de uma árvore de sintaxe para pelo menos uma sentença. Como exemplo de gramática ambígua, temos

$$E ::= E + E \mid E * E \mid \text{Numero}$$

onde Numero representa qualquer número natural. A sentença

$$3 + 4 * 5$$

pode ser gerada de duas formas usando derivação à esquerda:

$$\begin{aligned} E &\Longrightarrow E + E \Longrightarrow \text{Numero} + E \Longrightarrow 3 + E \Longrightarrow 3 + E * E \Longrightarrow 3 + \text{Numero} * E \\ &\Longrightarrow 3 + 4 * E \Longrightarrow 3 + 4 * \text{Numero} \Longrightarrow 3 + 4 * 5 \end{aligned}$$

e

$$\begin{aligned} E &\Longrightarrow E * E \Longrightarrow E + E * E \Longrightarrow \text{Numero} + E * E \Longrightarrow 3 + E * E \\ &\Longrightarrow 3 + \text{Numero} * E \Longrightarrow 3 + 4 * E \Longrightarrow 3 + 4 * \text{Numero} \Longrightarrow 3 + 4 * 5 \end{aligned}$$

Admitimos que +, * e os números naturais são terminais.

Ambigüidade é indesejável porque ela associa dois significados diferentes a uma sentença como “3 + 4 * 5”. No primeiro caso, a primeira derivação é

$$E \Longrightarrow E + E$$

o que implica que a multiplicação 4 * 5 será gerada pelo segundo E de “E + E”. Isto significa que a multiplicação deverá ser avaliada antes da soma, já que o resultado de 4 * 5 é necessário para E + E. Então, esta derivação considera que * possui maior precedência do que +, associando 23 como resultado de 3 + 4 * 5.

No segundo caso, a primeira derivação é

$$E \Longrightarrow E * E$$

e o primeiro E de “E * E” deve gerar 3 + 4, pois + vem antes de * na sentença “3 + 4 * 5”.² Sendo assim, esta seqüência de derivações admite que a soma deve ser calculada antes da multiplicação. Portanto é considerado que + possui maior precedência do que *.

2.3 Associatividade e Precedência de Operadores

Na linguagem S2 e na maioria das linguagens de programação, os operadores aritméticos +, -, * e / são associativos à esquerda. Isto é, “3 + 4 - 5” é avaliado como “(3 + 4) - 5”. Quando um operando, como 4, estiver entre dois operadores de mesma precedência (neste caso, + e -), ele será associado

²Se o segundo E gerasse um sinal +, teríamos algo como E * E + E, que não se encaixa na sentença 3 + 4 * 5.

ao operador mais à esquerda — neste caso, +. Por isto dizemos que os operadores aritméticos são associados à esquerda.

A gramática

$$\text{Expr} ::= \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term}$$

$$\text{Term} ::= \text{Term} * \text{Numero} \mid \text{Term} / \text{Numero} \mid \text{Numero}$$

com terminais +, -, *, / e Numero, associa todos os operadores à esquerda. Para entender como + é associativo à esquerda, basta examinar a produção

$$\text{Expr} ::= \text{Expr} + \text{Term}$$

Tanto Expr quanto Term podem produzir outras expressões. Mas Term nunca irá produzir um terminal +, como pode ser facilmente comprovado examinando-se a gramática. Então, todos os símbolos + obtidos pela derivação de “Expr + Term” se originarão de Expr. Isto implica que, em uma derivação para obter uma sentença com mais de um terminal +, como

$$2 + 6 + 1$$

o não-terminal Expr de “Expr + Term”, deverá se expandir para resultar em todos os símbolos +, exceto o último:

$$\text{Expr} \implies \text{Expr} + \text{Term} \implies \text{Expr} + 1 \implies \text{Expr} + \text{Term} + 1 \xRightarrow{+} 2 + 7 + 1$$

Implícito nesta derivação está que, para fazer a soma

$$\text{Expr} + \text{Term}$$

primeiro devemos fazer todas as somas em Expr, que está à direita de + em “Expr + Term”. Então, os operadores + à esquerda devem ser utilizados primeiro do que os + da direita, resultando então em associatividade à esquerda.

Voltando à gramática, temos que a regra

$$\text{Expr} ::= \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term}$$

produz, no caso geral, uma seqüência de somas e subtrações de não-terminais Term:

$$\text{Expr} \xRightarrow{+} \text{Term} + \text{Term} - \text{Term} + \text{Term} - \text{Term}$$

Então, as somas e subtrações só ocorrerão quando todos os não-terminais Term forem avaliados.

Examinando a regra para Term,

$$\text{Term} ::= \text{Term} * \text{Numero} \mid \text{Term} / \text{Numero} \mid \text{Numero}$$

descobrimos que Term nunca gerará um + ou -. Então, para avaliar Expr, devemos avaliar uma seqüência de somas e subtrações de Term. E, antes disto, devemos executar todas as multiplicações e divisões geradas por Term. Isto implica que * e / possuem maior precedência do que + e -.

Um operador * possuirá maior precedência do que + quando * tomar os seus operandos antes do que +. Isto é, se houver um operando (número) entre um * e um +, este será ligado primeiro ao *.

Um exemplo é a expressão

$$3 + 4 * 5$$

onde 4 está entre + e * e é avaliada como

$$3 + (4 * 5)$$

2.4 Modificando uma Gramática para Análise

As seções seguintes descrevem como fazer a análise sintática para uma dada gramática. O método de análise utilizado requer que a gramática :

1. não seja ambígua;
2. esteja fatorada à esquerda e;
3. não tenha recursão à esquerda.

O item 1 já foi estudado nas seções anteriores. O item 2, fatoração à esquerda, requer que a gramática não possua produções como

$$A ::= \beta \alpha_1 \mid \beta \alpha_2$$

onde β , α_1 e α_2 são seqüências de terminais e não-terminais. Uma gramática com produções

$$A ::= \beta_1 \mid \beta_1 \alpha_1$$

$$A ::= \beta_2 \mid \beta_2 \alpha_2$$

pode ser fatorada à esquerda transformando estas produções em

$$A ::= \beta_1 X_1 \mid \beta_2 X_2$$

$$X_1 ::= \epsilon \mid \alpha_1$$

$$X_2 ::= \epsilon \mid \alpha_2$$

onde ϵ corresponde à forma sentencial vazia. Observe que fatoração à esquerda é semelhante à fatoração matemática, como transformar $x + ax^2$ em $x(1 + ax)$.

O item 3 requer que a gramática não tenha recursão à esquerda. Uma gramática será recursiva à esquerda se ela tiver um não-terminal A de tal forma que exista uma derivação $A \xrightarrow{+} A\alpha$ para uma seqüência α de terminais e não-terminais.

Se uma gramática tiver produções

$$A ::= A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

onde nenhum β_i começa com A , a recursão à esquerda pode ser eliminada transformando-se estas produções em

$$A ::= \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' ::= \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

Observe que as primeiras produções geram uma seqüência de zero ou mais α_i iniciada por um único β_j , como

$$\beta_3 \alpha_3 \alpha_1 \alpha_5$$

$$\beta_5 \alpha_7 \alpha_1 \alpha_1 \alpha_1 \alpha_5 \alpha_3$$

$$\beta_2$$

Após a eliminação da recursão à esquerda, torna-se claro que um β_j aparecerá primeiro na forma sentencial derivada de A por causa das produções

$$A ::= \beta_j A'$$

e que haverá uma repetição de um ou mais α_i por causa das produções

$$A' ::= \alpha_i A'$$

A regra acima não funcionará quando A derivar $A\alpha$ em dois ou mais passos, como acontece na gramática

$$S ::= Aa \mid b$$

$$A ::= Ac \mid Sd \mid E$$

extraída de Aho, Sethi e Ullman [4]. S é recursivo à esquerda, pois $S \implies Aa \implies Sda$. Aho et al. apresentam um algoritmo (página 177) capaz de eliminar este tipo de recursão para algumas gramáticas. Este algoritmo não será estudado neste curso.

Utilizando as técnicas descritas acima, eliminaremos a recursão à esquerda e fatoraremos a gramática

$$\text{Expr} ::= \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term}$$

$$\text{Term} ::= \text{Term} * \text{Numero} \mid \text{Term} / \text{Numero} \mid \text{Numero}$$

Começamos eliminando a recursão à esquerda de Expr , considerando A igual a Expr , α_1 e α_2 iguais a “+ Term” e “- Term” e β_1 igual a Term . O resultado é

$$\text{Expr} ::= \text{Term Expr}'$$

$$\text{Expr}' ::= + \text{Term Expr}' \mid - \text{Term Expr}' \mid \epsilon$$

fazendo o mesmo com Term , obtemos

$$\text{Term} ::= \text{Numero Term}'$$

$$\text{Term}' ::= * \text{Numero Term}' \mid / \text{Numero Term}' \mid \epsilon$$

A gramática resultante já está fatorada.

Freqüentemente, a recursividade de um não terminal como Expr' é transformado em iteração pelo uso dos símbolos $\{ e \}$ na gramática, que indicam “repita zero ou mais vezes”. Assim,

$$A ::= a\{b\}$$

gera as sentenças

a
ab
abbbb

aplicando esta notação para a gramática anterior, obtemos

$$\begin{aligned} \text{Expr} &::= \text{Term} \{ (+|-) \text{Term} \} \\ \text{Term} &::= \text{Numero} \{ (*|/) \text{Numero} \} \end{aligned}$$

onde $+|-$ significa $+$ ou $-$.

2.5 Análise Sintática Descendente

A análise sintática (*parsing*) de uma seqüência de *tokens* determina se ela pode ou não ser gerada por uma gramática. A seqüência de *tokens* é gerada pelo analisador léxico a partir do programa fonte a ser compilado.

Existem diversos métodos de análise sintática e o método que estudaremos a seguir é chamado de descendente, pois ele parte do símbolo inicial da gramática e faz derivações buscando produzir como resultado final a seqüência de *tokens* produzida pelo analisador léxico. Um exemplo simples de análise descendente é analisar

$$2 + 3$$

a partir da gramática

$$\begin{aligned} E &::= T E' \\ E' &::= + T E' \mid \epsilon \\ T &::= N T' \\ T' &::= * N T' \mid \epsilon \\ N &::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

onde $+$, $*$ e os dígitos $0\dots9$ são terminais. Inicialmente, marcamos o primeiro símbolo da entrada com sublinhado e começamos a derivar a partir do símbolo inicial E :

$$\begin{array}{l} \underline{2} + 3\$ \\ E \end{array}$$

O $\$$ foi acrescentado ao final da entrada para representar o *token* “fim de arquivo”. Marcaremos com sublinhado o próximo símbolo sendo considerado da forma sentencial derivada de E . À esquerda do sublinhado estarão apenas terminais que já foram acasalados com a entrada.

Como só existe uma alternativa para derivar E , assim fazemos:

$$\begin{array}{l} \underline{2} + 3\$ \\ T E' \end{array}$$

Existe também uma única alternativa para T :

$$\begin{array}{l} \underline{2} + 3\$ \\ N T' E' \end{array}$$

Agora N possui 10 alternativas e escolheremos aquela que acasala (*match*) com o *token* corrente da entrada:

$$\begin{array}{l} \underline{2} + 3\$ \\ \underline{2} T' E' \end{array}$$

Quando o *token* corrente da entrada, também chamado de símbolo *lookahead*, for igual ao terminal da seqüência derivada de E , avançamos ambos os indicadores de uma posição, indicando que o *token* da

entrada foi aceito pela gramática:

$$2 \underline{+} 3\$$$

$$2 \underline{T'} E'\$$$

Agora, T' deveria derivar uma *string* começada por $+$. Como isto não é possível ($T' ::= * N T' \mid \epsilon$), escolhemos derivar T' para ϵ e deixar que $+$ seja gerado pelo próximo não terminal, E' :

$$2 \underline{+} 3\$$$

$$2 \underline{E'}$$

E' possui duas produções, $E' ::= + T E'$ e $E' ::= \epsilon$ e escolhemos a primeira por que ela acasala com a entrada:

$$2 \underline{+} 3\$$$

$$2 \underline{+} T E'$$

O $+$ é aceito resultando em

$$2 + \underline{3}\$$$

$$2 + \underline{T} E'$$

T é substituído pela sua única produção resultando em

$$2 + \underline{3}\$$$

$$2 + \underline{N} T' E'$$

Novamente, escolhemos a produção que acasala com a entrada, que é $N ::= 3$:

$$2 + \underline{3}\$$$

$$2 + \underline{3} T' E'$$

Aceitando 3 , temos

$$2 + \underline{3}\$$$

$$2 + 3 \underline{T'} E'$$

Como não existem mais caracteres a serem aceitos, ambos T' e E' derivam para ϵ :

$$2 + \underline{3}\$$$

$$2 + 3$$

Obtendo uma derivação a partir do símbolo inicial E igual à entrada estamos considerando que a análise sintática obteve sucesso.

Em alguns passos da derivação acima, o não terminal sendo considerado possuiu várias alternativas e escolhemos aquela que acasala a entrada. Por exemplo, em

$$\underline{2} + 3\$$$

$$\underline{N} T' E'$$

escolhemos derivar N utilizando $N ::= 2$ porque o *token* corrente da entrada era 2 .

No caso geral de análise, não é garantido que escolheremos a produção correta ($N ::= 2$ neste caso) a ser derivada, mesmo nos baseando no *token* corrente de entrada.

Como exemplo, considere a gramática

$$S ::= cAd$$

$$A ::= ab \mid a$$

tomada de um exemplo de Aho, Sethi e Ullman [4]. Com a entrada “cad” temos a seguinte análise:

$$\underline{c}ad\$$$

$$\underline{S}$$

$$\underline{c}ad\$$$

$$\underline{c}Ad$$

$$\underline{c}ad\$$$

$$\underline{c}Ad$$

```

cad$
cabd

```

```

cad$
cabd

```

Na última derivação há um erro de sintaxe: os dois terminais, da entrada e da forma sentencial derivada, são diferentes. Isto aconteceu porque escolhemos a produção errada de A . Portanto, temos que fazer um retrocesso (*backtracking*) à última derivação:

```

cad$
cAd

```

e escolher a segunda produção de A , que é $A ::= a$:

```

cad$
cad

```

que resulta em

```

cad$
cad

```

que aceita a entrada.

Em muitos casos, podemos escrever a gramática de tal forma que o retrocesso nunca seja necessário. Um analisador sintático descendente para este tipo de gramática é chamado de analisador preditivo. Em uma gramática para este tipo de analisador, dado que o símbolo corrente de entrada (*lookahead*) é “a” e o terminal a ver expandido (ou derivado) é A , existe uma única alternativa entre as produções de A ,

$$A ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

que começa com a *string* “a”. Se uma gramática não tiver nenhuma produção do tipo $A ::= \epsilon$, a condição acima será suficiente para garantir que um analisador preditivo pode ser construído para ela. As condições necessárias e suficientes para a construção de um analisador preditivo serão apresentadas no Capítulo 3.

As gramáticas de linguagens de programação em geral colocam uma palavra chave antes de cada comando exceto chamada de procedimento ou atribuição, facilitando a construção de analisadores preditivos. Utilizando-se uma gramática

```

Stmt ::= “if” Expr “then” Stmt “endif”
Stmt ::= “return” Expr
Stmt ::= “while” Expr “do” Stmt
Stmt ::= ident AssigOrCall
AssigOrCall ::= “=” Expr | “(” ExprList “)”

```

onde os símbolos entre aspas são terminais, pode-se analisar a entrada

```

if achou
then
  i = 1;
  return 0;
endif

```

sem retrocesso.

2.6 Análise Sintática Descendente Recursiva

Análise sintática descendente recursiva é um método de análise sintática descendente que emprega uma subrotina para cada não terminal da gramática. As subrotinas chamam-se entre si, em geral

havendo recursão entre elas.

Estudaremos apenas analisadores preditivos, isto é, sem retrocesso. Analisadores deste tipo podem ser construídos para, provavelmente, todas as linguagens de programação. Assim, não é uma limitação nos restringir a analisadores preditivos.

Uma gramática adequada para a construção de um analisador recursivo descendente não pode ser ambígua, deve estar fatorada à esquerda e não deve ter recursão à esquerda.

Estas condições são necessárias mas não suficientes para que possamos construir um analisador preditivo para a gramática. As condições necessárias e suficientes para que um analisador preditivo possa ser construído serão apresentadas nas próximas seções.

A gramática

$$\begin{aligned} \text{Expr} &::= \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term} \\ \text{Term} &::= \text{Term} * \text{Numero} \mid \text{Term} / \text{Numero} \mid \text{Numero} \\ \text{Numero} &::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

não é adequada para análise preditiva pois não está fatorada à esquerda e possui recursão à esquerda. Esta gramática foi modificada para

$$\begin{aligned} \text{Expr} &::= \text{Term Expr2} \\ \text{Expr2} &::= + \text{Term Expr2} \mid - \text{Term Expr2} \mid \epsilon \\ \text{Term} &::= \text{Numero Term2} \\ \text{Term2} &::= * \text{Numero Term2} \mid / \text{Numero Term2} \mid \epsilon \\ \text{Numero} &::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

na Seção 2.4, eliminando estes problemas.

O analisador preditivo para esta gramática é apresentada abaixo. Assume-se que exista um analisador léxico `nextToken` que coloca o token corrente em `token`.

Para cada não terminal da gramática existe uma função que o analisa. Sempre que o token corrente da entrada for igual ao terminal esperado pela função, a função `nextToken` é chamada para avançar para o próximo `token`. Isto é, quando o analisador estiver na configuração

$$\begin{array}{l} 2 \ \underline{\pm} \ 3\$ \\ 2 \ \underline{\pm} \ \text{Term Expr}' \end{array}$$

ele fará o teste

```
if ( token == mais_smb ) {
    nextToken();
    ...
}
```

aceitando o + e passando para o próximo `token` da entrada.

```
void expr()
{
    term();
    expr2();    // utilizamos expr2 para Expr'
}
```

```
void expr2()
{
    if ( token == mais_simb ) {
        // Expr' ::= + Term Expr'
        nextToken();
        term();
    }
}
```

```
    expr2();
  }
else
  if (token == menos_simb) {
    // Expr' ::= - Term Expr'
    nextToken();
    term();
    expr2();
  }
  /* se nenhum dos dois testes for verdadeiro (+ ou -), a alternativa
  escolhida sera Expr' ::= vazio e esta funcao nao deve fazer nada
  */
}

void term()
{
  if ( token == numero_smb ) {
    nextToken();
    term2();
  }
  else
    // Como nao ha alternativa Term ::= vazio, houve um erro
    erro();
}

void term2()
{
  switch (token) {
    case mult_smb:
      nextToken();
      if (token == numero_smb)
        nextToken();
      else
        erro();
      term2();
      break;
    case div_smb:
      nextToken();
      if (token == numero_smb)
        nextToken();
      else
        erro();
      term2();
      break;
    default:
      // Ok, existe producao Term' ::= vazio
  }
}
```

```
void erro()
{
    puts("Erro de sintaxe");
    exit(1);
}
```

Chapter 3

Análise Sintática Descendente Não Recursiva

3.1 Introdução

Um analisador descendente não recursivo utiliza uma pilha e uma tabela durante a análise.¹ A tabela é preenchida antes da análise de acordo com um método que será estudado posteriormente. A pilha contém, durante a análise, símbolos terminais e não terminais.

Um único analisador, mostrado na Figura 3.1, é utilizado para todas as linguagens, embora a tabela seja dependente da gramática. Dada a gramática

1. $E ::= T E'$
2. $E' ::= + T E'$
3. $E' ::= \epsilon$
4. $T ::= F T'$
5. $T' ::= * F T'$
6. $T' ::= \epsilon$
7. $F ::= (E)$
8. $F ::= \text{id}$

onde $+$, $*$, $($, $)$ e id são terminais, a tabela associada, que chamaremos de M , é

	id	()	+	*	eof
E	1	1	–	–	–	–
E'	–	–	3	2	–	3
T	4	4	–	–	–	–
T'	–	–	6	6	5	6
F	8	7	–	–	–	–

O símbolo – deve ser lido 0 (zero).

O algoritmo da Figura 3.1 será utilizado a seguir para analisar a sentença

$\text{id} + (\text{id} * \text{id})$

A cada passo da análise, mostraremos a pilha, os elementos da entrada ainda não analisados e a ação a ser tomada. Na pilha, o topo é colocado mais à esquerda. Assim, T estará no topo se a pilha for

E' T

Após uma operação desempilha, teremos

E'

¹Parte do texto e os exemplos deste capítulo foram retirados de [6] e [4].


```
void analiseNaoRecursiva()
{
    /* Analisa um programa cujos tokens sao fornecidos pelo objeto
       lex. A tabela de analise é a matriz M, ambos globais. */

    Pilha pilha;

    pilha.crie();
    pilha.empilhe(S); // S é o simbolo inicial

    while ( ! pilha.vazia() )
        if ( pilha.getTopo() == lex->token ) {
            pilha.desempilha();
            lex->nextToken();
        }
        else
            if ( M[ pilha.getTopo(), lex->token ] == 0 )
                ce->signal (erro_de_sintaxe_err);
            else {
                P = pilha.desempilha();
                empilhe todos os simbolos do lado direito da producao
                M[P, lex->token], da direita para a esquerda
            }

    if ( lex->token != eof_smb )
        ce->signal ( erro_de_sintaxe_err );
    else
        aceite a entrada

} // analiseNaoRecursiva
```

Figure 3.1: Algoritmo de análise não recursiva

Ao empilhar o lado direito de uma regra, como $F T'$ da regra $T ::= F T'$, primeiro invertemos os símbolos, obtendo $T' F$, e então acrescentaremos o resultado à pilha, resultando em $E' T' F$.

Com isto, simulamos a derivação

$$T E' \Longrightarrow F T' E'$$

Primeiro o símbolo mais à esquerda na sentença a ser derivada, T , é removido da pilha e então substituído pelo lado esquerdo da regra $T ::= F T'$.

A sentença $T E'$ é representada na pilha como $E' T$. Com a derivação $T E' \Longrightarrow F T' E'$, o topo T é desempilhado e os símbolos $T' F$ são empilhados, nesta ordem.

Estudaremos agora a análise da sentença

$$\text{id} + (\text{id} * \text{id})$$

utilizando o algoritmo da Figura 3.1. Cada passo da análise corresponde à aplicação de uma das regras descritas a seguir.

- emp r_i

Para desempilhar o topo da pilha e empilhar o lado direito da regra i . Esta regra é $M[\text{pilha.getTopo()}, \text{lex} \rightarrow \text{token}]$.

- desemp

Para desempilhar o topo da pilha e passar para o próximo símbolo. Esta ação será tomada quando o topo da pilha for um terminal igual ao *token* corrente da entrada.

- aceita

Para considerar a sentença válida. Este passo será utilizado quando a entrada for eof e a sentença a ser derivada for ϵ .

- erro

Quando o topo da pilha for terminal diferente do *token* corrente da entrada.

A análise da sentença $\text{id} + (\text{id} * \text{id})$ é detalhada na Figura 3.2.

Análise sintática descendente não recursiva pode ser utilizada com gramáticas que obedecem às mesmas restrições da análise recursiva: não ser ambígua, estar fatorada à esquerda e não ter recursão à esquerda.

3.2 A Construção da Tabela M

Para construir a tabela M utilizamos as relações **first** e **follow**, descritas a seguir.

$\text{first}(\alpha)$ é o conjunto dos terminais que iniciam α em uma derivação. Assim, utilizando a gramática

$$E ::= T E'$$

$$E' ::= + T E' \mid \epsilon$$

$$T ::= N T'$$

$$T' ::= * N T' \mid \epsilon$$

temos que

$$\text{first}(N) = \{ N \}$$

$$\text{first}(T') = \{ *, \epsilon \}$$

Pilha	Entrada	Ação
E	<u>id</u> + (id*id)	emp r ₁
E' T	<u>id</u> + (id*id)	emp r ₄
E' T' F	<u>id</u> + (id*id)	emp r ₈
E' T' id	<u>id</u> + (id*id)	desemp
E' T'	<u>±</u> (id*id)	emp r ₆
E'	<u>±</u> (id*id)	emp r ₂
E' T+	<u>±</u> (id*id)	desemp
E' T	(<u>id</u> *id)	emp r ₄
E' T' F	(<u>id</u> *id)	emp r ₇
E' T') E ((<u>id</u> *id)	desemp
E' T') E	<u>id</u> *id)	emp r ₁
E' T') E' T	<u>id</u> *id)	emp r ₄
E' T') E' T' F	<u>id</u> *id)	emp r ₈
E' T') E' T' id	<u>id</u> *id)	desemp
E' T') E' T'	<u>*</u> id)	emp r ₅
E' T') E' T' F *	<u>*</u> id)	desemp
E' T') E' T' F	<u>id</u>)	emp r ₈
E' T') E' T' id	<u>id</u>)	desemp
E' T') E' T')	emp r ₆
E' T') E')	emp r ₃
E' T'))	desemp
E' T'	eof	emp r ₆
E'	eof	emp r ₃
ε	eof	aceita

Figure 3.2: A análise da sentença $id + (id*id)$

$$\begin{aligned}\text{first}(T) &= \{ N \} \\ \text{first}(E') &= \{ +, \epsilon \} \\ \text{first}(E) &= \{ N \}\end{aligned}$$

A função **first** é definida formalmente como:

$$\begin{aligned}\text{first}: (N \cup \Sigma)^* &\longrightarrow \mathcal{P}(\Sigma \cup \{\epsilon\}) \\ \text{first}(\alpha) &= \{ a \mid \alpha \xrightarrow{*} a\beta, a \in \Sigma \cup \{\epsilon\} \}\end{aligned}$$

$\mathcal{P}(A)$ é o conjunto de todas as partes do conjunto A . Se $A = \{a, b\}$, $\mathcal{P}(A) = \{\emptyset, \{a\}, \{b\}, \{a, b\}\}$. Se A é um conjunto, A^* é o conjunto de todas as strings, cadeias ou sentenças² compostas por símbolos de A . Por exemplo, se $A = \{a, b\}$, $A^* = \{a, aa, aba, babb, \dots\}$. O símbolo $(N \cup \Sigma)^*$ significa o conjunto de todas as cadeias composta por símbolos terminais (Σ) e não terminais (N).

$\text{follow}(X)$ é o conjunto de terminais que sucedem X em qualquer derivação partindo do símbolo inicial S .

Se existir uma produção

$$Y ::= \alpha X a$$

então $a \in \text{follow}(X)$. O mesmo será válido se tivermos

$$Y ::= \alpha X Z$$

$$Z ::= aH \mid bG$$

e neste caso também teremos $b \in \text{follow}(X)$.

Formalmente, **follow** é definida como

$$\begin{aligned}\text{follow}: (N \cup \Sigma) &\longrightarrow \mathcal{P}(\Sigma \cup \{\text{eof}\}) \\ \text{follow}(X) &= \{ a \in \Sigma \cup \{\text{eof}\} \mid S \xrightarrow{*} \alpha X a \beta \text{ com } \alpha, \beta \in (N \cup \Sigma)^* \}\end{aligned}$$

first(X) é computado como:

1. se X for um terminal, $\text{first}(X) = \{ X \}$;
2. se $X ::= \epsilon$ for uma produção, adicione ϵ a $\text{first}(X)$;
3. se X for não terminal e $X ::= Y_1 Y_2 \dots Y_k$ for uma produção, então coloque f em $\text{first}(X)$ se, para algum i , $f \in \text{first}(Y_i)$ e $\epsilon \in \text{first}(Y_j)$, $j = 1, 2, \dots, i-1$. Se $\epsilon \in \text{first}(Y_j)$, $j = 1, 2, \dots, k$, então coloque ϵ em $\text{first}(X)$. Este item pode ser explicado como se segue.

Se Y_1 não derivar ϵ , então adicione $\text{first}(Y_1)$ a $\text{first}(X)$. Se $Y_1 \implies \epsilon$, adicione $\text{first}(Y_2)$ a $\text{first}(X)$. Se $Y_1 \implies \epsilon$ e $Y_2 \implies \epsilon$, adicione $\text{first}(Y_3)$ a $\text{first}(X)$ e assim por diante.

follow(X) é computado como:

1. coloque eof em $\text{follow}(S)$, onde S é o símbolo inicial;
2. se houver uma produção da forma $A ::= \alpha X \beta$, então $\text{first}(\beta)$, exceto ϵ , é colocado em $\text{follow}(X)$;
3. se houver uma produção da forma $A ::= \alpha X$ ou uma produção $A ::= \alpha X \beta$ e $\epsilon \in \text{first}(\beta)$, então colocaremos os elementos de $\text{follow}(A)$ em $\text{follow}(X)$.

A relação **first** deverá ser calculada antes de **follow** (por quê?). Para encontrar **first**(X) para todo $X \in N \cup \Sigma$ (todos os símbolos da gramática), aplique as regras acima até que nenhum terminal ou ϵ possa ser adicionado a nenhum conjunto **first**(Y). Isto é necessário porque a adição de elementos a um conjunto pode implicar, pela regra 3 de **first**, na adição de elementos a outros conjuntos.

²Aqui utilizados como sinônimos.

Para calcular $\text{follow}(A)$ para $A \in N \cup \{\text{eof}\}$ (não terminais e eof), aplique as regras anteriores até que nenhum elemento possa ser adicionado em qualquer conjunto $\text{follow}(B)$.

Calcularemos as relações first e follow para a gramática apresentada na introdução deste capítulo, que é

1. $E ::= T E'$
2. $E' ::= + T E'$
3. $E' ::= \epsilon$
4. $T ::= F T'$
5. $T' ::= * F T'$
6. $T' ::= \epsilon$
7. $F ::= (E)$
8. $F ::= \text{id}$

Os símbolos $+$, $*$, $($, $)$ e id são terminais. Observe que esta gramática obedece as restrições para a análise não recursiva, que são: não ser ambígua, estar fatorada à esquerda e não ter recursão à esquerda.

A relação first é:

$$\begin{aligned} \text{first}(E) &= \{ \text{id}, (\} \\ \text{first}(E') &= \{ +, \epsilon \} \\ \text{first}(T) &= \{ \text{id}, (\} \quad \text{Note que não calculamos a função } \text{first} \text{ completamente. Esta} \\ \text{first}(T') &= \{ *, \epsilon \} \\ \text{first}(F) &= \{ (, \text{id} \} \end{aligned}$$

função se aplica a qualquer cadeia de $(N \cup \Sigma)^*$. Então podemos calcular $\text{first}(T E')$ ou $\text{first}(* F T')$. No primeiro caso, $\text{first}(T E') = \text{first}(T)$, pois $\epsilon \notin \text{first}(T)$. No segundo caso, temos $\text{first}(* F T') = *$, pois $*$ é um terminal.

Os últimos conjuntos são calculados antes dos outros porque deles dependem os primeiros: $\text{first}(A)$ dependerá de $\text{first}(B)$ se houver uma produção $A ::= B\alpha$.

A relação follow é:

$$\begin{aligned} \text{follow}(E) &= \{ \text{eof},) \} \\ \text{follow}(E') &= \{ \text{eof},) \} \\ \text{follow}(T) &= \{ +, \text{eof},) \} \\ \text{follow}(T') &= \{ +, \text{eof},) \} \\ \text{follow}(F) &= \{ *, +, \text{eof},) \} \end{aligned}$$

A relação follow para os terminais não foi calculada, apesar de ser válida, porque não será utilizada no exemplo a seguir.

A tabela M utilizada pelo algoritmo de análise sintática não recursiva é calculada como se segue. Para cada produção da forma $A ::= \alpha$, faça:

1. para cada $a \in \text{first}(\alpha)$, adicione o número da regra
 $A ::= \alpha$
em $M[A, a]$;
2. Se $\epsilon \in \text{first}(\alpha)$, adicione o número da regra
 $A ::= \alpha$
em $M[A, b]$ para cada $b \in \text{follow}(A)$;

Torne indefinidas todas as outras entradas:

$$M[A, a] = 0$$

A tabela construída de acordo com estas regras foi apresentada na introdução deste capítulo.

3.3 Gramáticas LL(1)

Seja G uma gramática sem recursão à esquerda e fatorada à esquerda a partir da qual foi construída uma tabela M pelo método apresentado na seção anterior. Se G for ambígua, haverá uma entrada $M[X, a]$ com mais de um elemento. Por exemplo, considere a gramática

1. $S ::= A$
2. $S ::= B$
3. $A ::= a C$
4. $B ::= ab$
5. $C ::= b$

Podemos obter a sentença “ab” com duas derivações à esquerda diferentes:

$$S \Rightarrow A \Rightarrow a C \Rightarrow a b$$

$$S \Rightarrow B \Rightarrow a b$$

Isto caracteriza a ambigüidade. A tabela de análise desta gramática é

	a	b	eof
S	1, 2		
A	3		
B	4		
C		5	

Como $M[S, a]$ é $\{ 1, 2 \}$, então poderemos utilizar a produção 1 ou 2 para expandir S quando o símbolo corrente da entrada for “a”.

Contudo, uma gramática pode ter mais de um número em certa entrada da tabela e não ser ambígua. Por exemplo, a gramática

1. $S ::= A$
2. $S ::= B$
3. $A ::= a C$
4. $B ::= a$
5. $C ::= b$

possui a mesma tabela de análise que a gramática anterior, mas não é ambígua. É fácil ver porquê. A linguagem desta gramática só possui “ab” e “a” como sentenças (a gramática só produz estas sentenças) e só existe uma forma de produzir cada uma delas:

$$S \Rightarrow A \Rightarrow a C \Rightarrow a b$$

$$S \Rightarrow B \Rightarrow a$$

A gramática

1. $S ::= \text{“if” } E \text{ “then” } S S'$
2. $S ::= a$

3. $S' ::= \text{“else” } S$
4. $S' ::= \epsilon$
5. $E ::= b$

tomada de Aho et al. [4] possui a seguinte tabela de análise:

	a	b	“else”	“if”	“then”	eof
S	2			1		
S'			3, 4			4
E		5				

A entrada $M[S', \text{“else”}]$ possui dois elementos, implicando em considerar o “else” da sentença `if b then a if b then a else a` associado ao último `if` (3) ou ao primeiro (4).

Uma gramática será chamada de LL(1) se a sua tabela de análise tiver no máximo um elemento por cada entrada $M[X, a]$. O primeiro L de “LL(1)” informa que a análise para a gramática é feita da esquerda (*Left*) para a direita. O segundo L implica que, na análise, é produzida a derivação mais à esquerda sendo que em cada passo da derivação a ação a ser tomada é decidida com base em um (“1” de “LL(1)”) símbolo da entrada (o *lookahead*).

Pode ser provado que uma gramática G é LL(1) se e somente se sempre que $A ::= \alpha \mid \beta$ forem duas produções distintas de G , os itens a seguir são verdadeiros.³

1. Não há um terminal a pertencente ao mesmo tempo a $\text{first}(\alpha)$ e $\text{first}(\beta)$.
2. Ou $\alpha \xRightarrow{*} \epsilon$ ou $\beta \xRightarrow{*} \epsilon$, nunca estas duas condições ao mesmo tempo.
3. Se $\beta \xRightarrow{*} \epsilon$, então α não derivará qualquer *string* começando com um terminal que pertence a $\text{follow}(A)$.

Discutiremos porque a desobediência a um destes itens implica na ambigüidade da gramática.

1. Se $\alpha \xRightarrow{*} a$ e $\beta \xRightarrow{*} a$ puderem ocorrer e tivermos que expandir A com a como o símbolo corrente da entrada (*lookahead*), teremos duas produções a escolher: $A ::= \alpha$ e $A ::= \beta$.
2. Se for possível $\alpha \xRightarrow{*} \epsilon$ e $\beta \xRightarrow{*} \epsilon$ e tivermos que expandir A com b como símbolo corrente de entrada, $b \notin \text{first}(A)$, teremos novamente duas produções a escolher.
3. Se β puder derivar ϵ ($\beta \xRightarrow{*} \epsilon$), $b \in \text{follow}(A)$ e α puder derivar b ($\alpha \xRightarrow{*} b\gamma$), então tanto $A ::= \alpha$ quanto $A ::= \beta$ poderão ser escolhidas para expandir A quando b for o símbolo corrente da entrada. É claro que a utilização de $A ::= \alpha$ seria correta, já que $b \in \text{first}(\alpha)$. Para enxergar porque a utilização de $A ::= \beta$ produziria o resultado correto, considere o seguinte exemplo: $\gamma A \phi$ está sendo expandido para reconhecer a entrada, b é o token corrente, γ já foi reconhecido, A deve ser derivado, $b \in \text{first}(\phi)$ e, portanto, $b \in \text{follow}(A)$. Escolhendo $A ::= \beta$ e $\beta \xRightarrow{*} \epsilon$, teríamos $\gamma A \phi \xRightarrow{+} \gamma \phi$ onde b seria reconhecido por ϕ .

Nem todas as gramáticas são ou podem ser transformadas em LL(1). Este é o principal obstáculo ao uso de analisadores preditivos. Felizmente, o problema de entradas na tabela com mais de um elemento pode ser resolvido escolhendo-se um deles, removendo assim a ambigüidade.

³Todo este trecho é uma tradução quase literal de um parágrafo da página 192 de Aho et al. [4].

Analisadores ascendentes (*botton-up*), que derivam o símbolo inicial a partir dos terminais, podem ser utilizados com um número maior de gramáticas do que analisadores descendentes preditivos. Estes analisadores são, em geral, gerados automaticamente por *geradores de analisadores sintáticos* a partir da gramática da linguagem. Analisadores ascendentes empregam grandes tabelas que seriam difíceis de serem gerados por pessoas sem o auxílio de programas. Como exemplos de geradores de analisadores sintáticos, temos o CUP, JavaCC, YACC, Bison, ANTLR.

Bibliography

- [1] Stroustrup, Bjarne. *The C++ Programming Language*. Second Edition, Addison-Wesley, 1991.
- [2] Lippman, Stanley B. *C++ Primer*. Addison-Wesley, 1991.
- [3] Deitel, H.M. e Deitel P.J. *C++ How to Program*. Prentice-Hall, 1994.
- [4] Aho, Alfred e Sethi, Rave e Ullman, Jeffrey. *Compilers: Principles, Techniques and Tools*. Addison-Wesley Publishing Company, 1986.
- [5] Gamma, Erich; Helm, Richard; Johnson, Ralph e Vlissides, John. *Design Patterns: Elements of Reusable Object-Oriented Software*. Professional Computing Series, Addison-Wesley, Reading, MA, 1994.
- [6] Zorzo, Sérgio Donizetti. Notas de aula de Construção de Compiladores, 1995.
- [7] Pittman, Thomas; Peter, James. *The Art of Compiler Design: Theory and Practice*. Prentice Hall, 1992.
- [8] Guimarães, José de Oliveira. *The Green Language*. Disponível em <http://www.dc.ufscar.br/jose/green/green.htm>.