

Alan D. Barroso, Kenji Sakata Jr

Anlisador Sintático

Universidade de São Paulo
Escola Politécnica

Linguagens e Compiladores
PCS2056

São Paulo
2013

1 Introdução

O analisador sintático é o responsável por reconhecer a estrutura da linguagem, formada a partir da gramática. Sua ação começa uma vez que o módulo léxico já reconheceu todas as partículas atômicas que compõem o código-fonte.

Na realidade, é o analisador sintático quem invoca a ação do módulo léxico em um primeiro momento a fim de identificar as partículas do texto. Em seguida, ele, o analisador sintático, faz uma segunda leitura do código-fonte, mas, desta vez, ele tenta descrever como essas partículas estão estruturadas e organizadas no texto. A estrutura que é composta por esta segunda leitura é chamada de árvore de sintaxe e consiste unicamente em uma forma alternativa de representar o código-fonte, mas a qual é compreensível pelo compilador.

Em seguida, quando da análise semântica, é novamente o módulo sintático quem será o responsável por invocar o analisador semântico, o qual traduzirá a árvore de sintaxe em ações concretas, executáveis pelo computador.

Outras funções tipicamente atribuídas ao analisador sintático são detecção de erros de sintaxe, recuperação de erros, correção de erros, ativação de rotinas de síntese do código-objeto, entre outras.

Existem outras formas de delegar essas tarefas aos três módulos canônicos de um compilador clássico, mas podemos concluir que, por conta de nossas decisões de projeto, o compilador em desenvolvimento é orientado à sintaxe. Isso porque o módulo sintático do compilador é o pivô central que coordena o fluxo sequencial de ações que permitem traduzir uma linguagem de alto nível em linguagem de máquina. É ele quem determina qual é a etapa do processo de compilação que está sendo executada e quem é o responsável por ela.

Assim, após definir a linguagem e desenvolver o analisador léxico, tratamos de construir o módulo de análise sintática, cujo método de reconhecimento foi baseado em autômato de pilha estruturado (APE).

2 Geração Automática dos Autômatos

2.1 Notação de Wirth Reduzida

Na primeira etapa para a criação do APE, foram utilizadas as expressões fundamentais criadas a partir da descrição reduzida em notação de Wirth. A definição das expressões seguiu o critério dos não-terminais com recursividade central, sendo possível concluir que as máquinas finais tratariam a sequência de comandos, as expressões e o programa principal.

Assim, a linguagem representada pela notação de Wirth reduzida pode ser encontrada no arquivo "wirth_reduzido.txt" que acompanha esse documento.

2.2 Lista de submáquinas do APE

2.2.1 Lista de transições

A partir da linguagem representada em notação de Wirth reduzida, utilizamos o programa do site indicado ¹ para gerar automaticamente a tabela de transições as quais as submáquinas deveriam executar de forma a caracterizar a linguagem.

O programa gera diversos tipos de saídas. Cuida lembrar que somente consideramos a saída que é uma tabela de transição reduzida, de forma a otimizar o processamento do compilador e facilitar a leitura e boa compreensão da representação dos autômatos, seja ela na forma tabular ou gráfica.

Grosso modo, o programa utiliza três etapas para reduzir a tabela de transições. A ordem em que elas devem ser executadas é a seguinte:

- Eliminação das transições em vazio;
- Eliminação dos estados não-acessíveis;
- Eliminação dos estados equivalentes.

A saída da tabela referente à submáquina programa:

¹ <http://radiant-fire-72.herokuapp.com/>

```
1 initial: 0
2 final: 29
3 (0, "const") -> 1
4 (0, "typedef") -> 2
5 (0, "struct") -> 3
6 (0, "vars") -> 4
7 (0, "function") -> 5
8 (0, "program") -> 6
9 (1, identificador) -> 7
10 (2, identificador) -> 23
11 (3, identificador) -> 30
12 (4, identificador) -> 53
13 (4, "bool") -> 53
14 (4, "int") -> 53
15 (4, "char") -> 53
16 (4, "float") -> 53
17 (5, identificador) -> 31
18 (5, "bool") -> 31
19 (5, "int") -> 31
20 (5, "char") -> 31
21 (5, "float") -> 31
22 (5, "void") -> 31
23 (6, "begin") -> 8
24 (7, "=") -> 9
25 (8, "end") -> 10
26 (8, "vars") -> 11
27 (8, comando) -> 12
28 (9, inteiro) -> 13
29 (9, float) -> 13
30 (9, "true") -> 13
31 (9, "false") -> 13
32 (9, char) -> 13
33 (10, "program") -> 29
34 (11, identificador) -> 15
35 (11, "bool") -> 15
36 (11, "int") -> 15
37 (11, "char") -> 15
38 (11, "float") -> 15
39 (12, "end") -> 10
40 (12, comando) -> 12
41 (13, ";") -> 14
42 (14, identificador) -> 7
43 (14, "end") -> 16
44 (15, identificador) -> 17
45 (16, "const") -> 18
46 (17, ";") -> 19
47 (17, "[") -> 20
48 (17, ",") -> 15
49 (18, "typedef") -> 2
50 (18, "struct") -> 3
51 (18, "vars") -> 4
```

```
52 (18, "function") -> 5
53 (18, "program") -> 6
54 (19, identificador) -> 15
55 (19, "end") -> 22
56 (19, "bool") -> 15
57 (19, "int") -> 15
58 (19, "char") -> 15
59 (19, "float") -> 15
60 (20, inteiro) -> 21
61 (21, "]" ) -> 17
62 (22, "vars") -> 12
63 (23, "=") -> 24
64 (24, identificador) -> 25
65 (24, "bool") -> 25
66 (24, "int") -> 25
67 (24, "char") -> 25
68 (24, "float") -> 25
69 (25, ";" ) -> 26
70 (26, identificador) -> 23
71 (26, "end") -> 27
72 (27, "typedef") -> 28
73 (28, "struct") -> 3
74 (28, "vars") -> 4
75 (28, "function") -> 5
76 (28, "program") -> 6
77 (30, "begin") -> 32
78 (31, identificador) -> 33
79 (32, identificador) -> 34
80 (32, "bool") -> 34
81 (32, "int") -> 34
82 (32, "char") -> 34
83 (32, "float") -> 34
84 (33, "(" ) -> 35
85 (34, identificador) -> 36
86 (35, identificador) -> 37
87 (35, "bool") -> 37
88 (35, "int") -> 37
89 (35, "char") -> 37
90 (35, "float") -> 37
91 (35, ")" ) -> 38
92 (36, ";" ) -> 39
93 (36, "[" ) -> 40
94 (36, "," ) -> 34
95 (37, identificador) -> 59
96 (37, "byref") -> 60
97 (38, "begin") -> 41
98 (39, identificador) -> 34
99 (39, "end") -> 51
100 (39, "bool") -> 34
101 (39, "int") -> 34
102 (39, "char") -> 34
```

```
103 (39, "float") -> 34
104 (40, inteiro) -> 45
105 (41, "end") -> 42
106 (41, "vars") -> 43
107 (41, comando) -> 44
108 (42, "function") -> 58
109 (43, identificador) -> 46
110 (43, "bool") -> 46
111 (43, "int") -> 46
112 (43, "char") -> 46
113 (43, "float") -> 46
114 (44, "end") -> 42
115 (44, comando) -> 44
116 (45, "]" ) -> 36
117 (46, identificador) -> 47
118 (47, ";" ) -> 48
119 (47, "[" ) -> 49
120 (47, "," ) -> 46
121 (48, identificador) -> 46
122 (48, "end") -> 52
123 (48, "bool") -> 46
124 (48, "int") -> 46
125 (48, "char") -> 46
126 (48, "float") -> 46
127 (49, inteiro) -> 50
128 (50, "]" ) -> 47
129 (51, "struct") -> 28
130 (52, "vars") -> 44
131 (53, identificador) -> 54
132 (54, ";" ) -> 55
133 (54, "[" ) -> 56
134 (54, "," ) -> 53
135 (55, identificador) -> 53
136 (55, "end") -> 63
137 (55, "bool") -> 53
138 (55, "int") -> 53
139 (55, "char") -> 53
140 (55, "float") -> 53
141 (56, inteiro) -> 57
142 (57, "]" ) -> 54
143 (58, "function") -> 5
144 (58, "program") -> 6
145 (59, "[" ) -> 61
146 (59, "," ) -> 62
147 (59, ")" ) -> 38
148 (60, identificador) -> 59
149 (61, inteiro) -> 64
150 (62, identificador) -> 37
151 (62, "bool") -> 37
152 (62, "int") -> 37
153 (62, "char") -> 37
```

```
154 (62, "float") -> 37
155 (63, "vars") -> 58
156 (64, "]" ) -> 59
```

../1-linguagem/notacoes/JFLAP/programa/programa.txt

A saída referente à submáquina comando:

```
1 initial: 0
2 final: 12
3 (0, identificador) -> 1
4 (0, "call") -> 2
5 (0, "if") -> 3
6 (0, "while") -> 4
7 (0, "scan") -> 5
8 (0, "print") -> 6
9 (0, "return") -> 7
10 (1, "[" ) -> 8
11 (1, "." ) -> 9
12 (1, "=" ) -> 10
13 (1, "+=" ) -> 10
14 (1, "-=" ) -> 10
15 (1, "*=" ) -> 10
16 (1, "/=" ) -> 10
17 (2, identificador) -> 18
18 (3, "(" ) -> 22
19 (4, "(" ) -> 28
20 (5, identificador) -> 14
21 (6, expressao) -> 13
22 (7, expressao) -> 11
23 (7, ";" ) -> 12
24 (8, expressao) -> 16
25 (9, identificador) -> 1
26 (10, expressao) -> 11
27 (11, ";" ) -> 12
28 (13, ";" ) -> 12
29 (13, "," ) -> 6
30 (14, "[" ) -> 15
31 (14, "." ) -> 5
32 (14, ";" ) -> 12
33 (14, "," ) -> 5
34 (15, expressao) -> 17
35 (16, "]" ) -> 1
36 (17, "]" ) -> 14
37 (18, "(" ) -> 19
38 (19, expressao) -> 20
39 (19, ")" ) -> 11
40 (20, "," ) -> 21
41 (20, ")" ) -> 11
42 (21, expressao) -> 20
43 (22, expressao) -> 23
44 (22, ")" ) -> 24
```

```

45 (23, ")") -> 24
46 (24, "then") -> 25
47 (25, comando) -> 25
48 (25, "else") -> 26
49 (25, "end") -> 27
50 (26, comando) -> 26
51 (26, "end") -> 27
52 (27, "if") -> 12
53 (28, expressao) -> 29
54 (28, ")") -> 30
55 (29, ")") -> 30
56 (30, "do") -> 31
57 (31, comando) -> 31
58 (31, "end") -> 32
59 (32, "while") -> 12

```

../1-linguagem/notacoes/JFLAP/comando/comando.txt

E, finalmente, a saída referente à submáquina expressão:

```

1 initial: 0
2 final: 3, 4, 9
3 (0, "!") -> 1
4 (0, "-") -> 1
5 (0, "(") -> 2
6 (0, inteiro) -> 3
7 (0, float) -> 3
8 (0, "true") -> 3
9 (0, "false") -> 3
10 (0, char) -> 3
11 (0, identificador) -> 4
12 (1, "(") -> 2
13 (1, inteiro) -> 3
14 (1, float) -> 3
15 (1, "true") -> 3
16 (1, "false") -> 3
17 (1, char) -> 3
18 (1, identificador) -> 4
19 (2, expressao) -> 8
20 (3, "-") -> 0
21 (3, "^") -> 0
22 (3, "*") -> 0
23 (3, "/") -> 0
24 (3, "+") -> 0
25 (3, "==") -> 0
26 (3, "!=") -> 0
27 (3, "<") -> 0
28 (3, ">") -> 0
29 (3, "<=") -> 0
30 (3, ">=") -> 0
31 (3, "and") -> 0
32 (3, "or") -> 0

```



```

33 (4, "-" ) -> 0
34 (4, "(" ) -> 5
35 (4, "[" ) -> 6
36 (4, "." ) -> 7
37 (4, "^" ) -> 0
38 (4, "*" ) -> 0
39 (4, "/" ) -> 0
40 (4, "+" ) -> 0
41 (4, "==" ) -> 0
42 (4, "!=" ) -> 0
43 (4, "<" ) -> 0
44 (4, ">" ) -> 0
45 (4, "<=" ) -> 0
46 (4, ">=" ) -> 0
47 (4, "and" ) -> 0
48 (4, "or" ) -> 0
49 (5, expressao) -> 11
50 (5, ")" ) -> 3
51 (6, expressao) -> 10
52 (7, identificador) -> 9
53 (8, ")" ) -> 3
54 (9, "-" ) -> 0
55 (9, "[" ) -> 6
56 (9, "." ) -> 7
57 (9, "^" ) -> 0
58 (9, "*" ) -> 0
59 (9, "/" ) -> 0
60 (9, "+" ) -> 0
61 (9, "==" ) -> 0
62 (9, "!=" ) -> 0
63 (9, "<" ) -> 0
64 (9, ">" ) -> 0
65 (9, "<=" ) -> 0
66 (9, ">=" ) -> 0
67 (9, "and" ) -> 0
68 (9, "or" ) -> 0
69 (10, "]" ) -> 9
70 (11, ")" ) -> 3
71 (11, "," ) -> 12
72 (12, expressao) -> 11

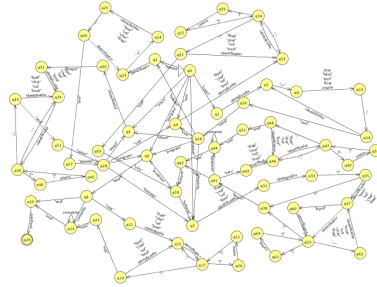
```

../1-linguagem/notacoes/JFLAP/expressao/expressao.txt

2.2.2 Lista dos autômatos

Nesta seção listaremos simplesmente os autômatos resultantes representados de uma maneira gráfica a partir das tabelas de transições da seção anterior.

Figura 1 – Autômato representando a submáquina programa



3 Comentários sobre a implementação

3.1 Desenvolvimento do módulo sintático

Partindo do modelo da notação de Wirth reduzida, nosso analisador sintático baseia-se na interação entre três máquinas de estados, representando as regras finais, que podem se chamar mutuamente, o que é previsto por estarmos trabalhando com um Autômato de Pilha Estruturado (APE).

O módulo sintático, portanto, começa suas atividades com a premissa que as três máquinas de estados já estejam criadas e inicializadas. Pressupõe-se que essa inicialização seja executada no programa principal (função `main`) do compilador. Essa inicialização ocorre de forma análoga ao que realizamos no módulo léxico e é baseada em tabelas de transições, as quais encontram-se descritas em arquivos externos, facilitando futuras modificações.

Utiliza-se a subrotina `verify_syntax()`, a qual já capaz de dizer se o compilador reconhece a linguagem. Mais detalhadamente, essa verificação tem como principal atividade o lançamento da função `recognize()`, utilizando a máquina que representa "programa" e o primeiro token lido pelo analisador léxico como parâmetros.

Grosso modo, o que esta função realiza é, dado o token passado via parâmetro, consulta-se sua tabela de transições internas e verifica-se se há alguma transição possível. Se sim, a transição é realizada, o estado do autômato é atualizado e um novo token é lido.

Cuida lembrar que há casos em que a transição chama uma nova submáquina, o que também é especificado pela tabela. Quando isso ocorre, a nova máquina é chamada através da mesma função `recognize`, mas passando a nova máquina como parâmetro e o mesmo token por referência. Em outros termos, não há necessidade de lookahead e, por serem armazenados em memória quando da chamada de submáquina, os tokens são consumidos sob demanda.

Por outro lado, se não houver transição possível na tabela para um dado token, a função verifica se a máquina encontra-se em estado de aceitação, retornando verdadeiro, ou não, retornando falso. O valor falso traduz o caso em que houve um erro sintático na estrutura do código-

fonte e o erro é rapidamente propagado para as máquinas de estados que a chamaram, parando a execução do módulo. Já o valor verdadeiro, faz com que a máquina que a chamou continue sua execução a partir do próximo estado previsto pela tabela em que o autômato estava a esperar (resincronização); se o autômato é a raiz ("programa") e já não há mais tokens a serem lidos, o compilador reconhece a sintaxe da linguagem.

3.2 Integração com o compilador

Devido à alta modularização configurada pela arquitetura escolhida, o programa principal que representa o compilador possui uma estrutura bastante simples.

Primeiramente, o módulo léxico é inicializado, passando o código-fonte a ser analisado como parâmetro. Em seguida, inicializa-se a Tabela de Símbolos e a Tabela de Palavras-chaves. Até aqui, não há diferenças em relação à entrega do analisador léxico.

A próxima ação, então, inicializa o módulo sintático, o que consiste basicamente na criação das máquinas de estados que representam as expressões de Wirth a partir dos arquivos externos, como descrito na seção anterior.

Finalmente, chama-se a subrotina `verify_syntax()`, a qual foi explicada anteriormente e é capaz de decidir se o código-fonte pertence ou não à linguagem definida pela gramática.