

TRABALHO 2 - ANÁLISE SINTÁTICA

Arthur Correia, Matheus Queiroz, Pedro Marcos Estrela

Dezembro, 2025

Introdução

Este projeto tem como objetivo o desenvolvimento de um analisador sintático para a linguagem de programação imperativa criada pelo grupo no Trabalho 1, utilizando as ferramentas **Flex** e **Bison**, no contexto da disciplina MATA61 – Compiladores.

Enquanto o analisador léxico (Trabalho 1) preocupava-se apenas em identificar os tokens individualmente, o analisador sintático aqui apresentado define a estrutura gramatical da linguagem. O objetivo final é receber o código-fonte como entrada e gerar uma **Árvore Sintática** que represente a hierarquia dos comandos, declarações e expressões.

A linguagem suporta tipagem estática, vetores, condicionais, laços e funções.

Descrição do Código e Implementação

O analisador sintático foi desenvolvido utilizando a ferramenta Bison, que trabalha em conjunto com o analisador léxico (Flex) produzido na etapa anterior. A estrutura do código é reconhecida e construída a partir dos tokens individuais (folhas) até formar a estrutura completa do programa (raiz).

A principal função deste analisador é validar a gramática da linguagem e construir a Árvore Sintática, que representa visualmente a hierarquia e a lógica do código fonte fornecido.

Estrutura de Dados da Árvore

Para permitir a visualização hierárquica exigida, foi definida uma estrutura de dados em C chamada `ASTNode`. Esta estrutura é responsável por armazenar as informações de cada nó da árvore:

- **Tipo:** A categoria do nó (ex: "declaracao", "if_statement", "operacao").
- **Valor:** O conteúdo textual, caso exista (ex: nome da variável "x", valor numérico "10").
- **Filhos:** Uma lista de ponteiros para outros nós, permitindo o aninhamento de comandos.

No Bison, a união de tipos (`%union`) foi configurada para manipular esses nós, permitindo que cada regra gramatical retorne um ponteiro para um nó da árvore.

```
1 %union {  
2     int intval;  
3     double floatval;  
4     char *strval;
```

```
5     struct ASTNode *node;
6 }
```

Dessa forma, o Flex identifica os tokens e envia os valores brutos (strings ou números), enquanto o Bison organiza esses valores criando e conectando os nós da árvore.

Regras de Produção e Justificativas

Estrutura do Programa

A regra inicial da gramática é o **programa**. Ela define que um código válido é composto por uma lista sequencial de comandos.

```
1 programa:
2     lista_comandos {
3         ASTNode *root = create_node("programa", NULL);
4         add_child(root, $1);
5         print_tree(root, 0);
6         free_tree(root);
7     }
8 ;
```

Lista de Comandos e Blocos

A construção do programa ocorre de forma recursiva através da regra **lista_comandos**, que permite o encadeamento infinito de instruções. Além disso, a regra **bloco** define o escopo de execução, fundamental para estruturas condicionais e laços.

```
1 lista_comandos:
2     lista_comandos comando {
3         $$ = $1;
4         add_child($$, $2); }
5     | comando {
6         $$ = create_node("lista_comandos", NULL);
7         add_child($$, $1); }
8     ;
9 bloco:
10    T_LCHAVE lista_comandos T_RCHAVE { $$ = $2; }
11    | T_LCHAVE T_RCHAVE { $$ = create_node("bloco_vazio", NULL); }
12    ;
```

Gerenciamento de Variáveis e Vetores

Além das variáveis escalares simples, a linguagem suporta vetores através da palavra-chave `serpente`.

```
1 declaracao_serpente:
2     T_SERPENTE_KW T_ID T_LPAREN tipo T_RPAREN T_ATRIBUICAO T_LPAREN
3     elementos_serpente T_RPAREN {
4         $$ = create_node("serpente_declaracao", $2);
5         add_child($$, $4);
6         add_child($$, $8);
7     }
8 ;
9 elementos_serpente:
10    elementos_serpente T_VIRGULA expressao { ... }
11    | expressao { ... }
12 ;
```

A gramática exige inicialização explícita para vetores, garantindo que a estrutura `serpente` seja sempre instanciada com seus elementos, facilitando a alocação de memória na geração da AST.

Estruturas de Repetição

A linguagem implementa dois tipos de laços: `enquanto` (baseado em condição) e `ate` (iterativo).

```
1 comando_enquanto:
2     T_ENQUANTO T_LPAREN expressao T_RPAREN bloco {
3         $$ = create_node("enquanto_loop", NULL);
4         ...
5     }
6 ;
7
8 comando_ate:
9     T_ATE T_LPAREN expressao T_RPAREN bloco {
10        $$ = create_node("ate_loop", NULL);
11        ...
12    }
13 ;
```

Ambas as estruturas segregam claramente a *condição* de parada do *bloco* de execução, permitindo que a AST represente o fluxo de controle de maneira estruturada.

Funções e Entrada/Saída

As funções são cidadãos de primeira classe na estrutura sintática, possuindo tipo de retorno e lista de parâmetros.

```
1 definicao_funcao:
2     tipo T_ID T_LPAREN lista_parametros T_RPAREN bloco {
3         $$ = create_node("funcao_definicao", $2);
4         add_child($$, $1);
5         add_child($$, $4);
6         add_child($$, $6);
7     }
8 ;
9
10 comando_devolva:
11     T_DEVOLVA expressao { ... } ;
12
13 comando_receba:
14     T RECEBA T_LPAREN T_ID T_RPAREN { ... } ;
15
16 comando_imprima:
17     T_IMPRIMA T_LPAREN expressao T_RPAREN { ... } ;
```

A regra `definicao_funcao` captura toda a assinatura do método antes de processar seu corpo. Os comandos de I/O (`receba`, `imprima`) e retorno (`devolva`) são tratados como instruções atômicas.

Expressões e Hierarquia de Operadores

Para garantir a precedência correta das operações matemáticas e lógicas sem o uso de parênteses excessivos, a gramática foi estratificada em três níveis: `expressao`, `termo` e `fator`.

```
1 expressao:
2     expressao_simples
3     | expressao_simples T_IGUAL expressao_simples { ... }
4     /* ... outros operadores relacionais ... */
5 ;
6
7 expressao_simples:
8     termo
9     | expressao_simples T_SOMA termo { ... }
10    | expressao_simples T_SUB termo { ... }
11    | expressao_simples T_OU termo { ... }
12 ;
13
14 termo:
15     fator
16     | termo T_MULT fator { ... }
```

```

17    | termo T_DIV fator { ... }
18    | termo T_E fator { ... }
19    ;
20
21 fator:
22     T_LPAREN expressao T_RPAREN { $$ = $2; }
23     | T_ID { ... }
24     | T_INTEIRO_LIT { ... }
25     /* ... outros literais ... */
26

```

Justificativa:

- **Fator:** É a unidade indivisível (números, variáveis) ou expressões priorizadas por parênteses.
- **Termo:** Processa multiplicações, divisões e o operador lógico 'E', que possuem precedência sobre soma/subtração.
- **Expressão Simples:** Processa somas, subtrações e o operador lógico 'OU'.
- **Expressão:** O nível mais alto, onde ocorrem as comparações relacionais (igualdade, maior que, etc.), que têm a menor precedência.

Essa cascata garante que a expressão $2 + 3 * 4$ seja parseada corretamente como $2 + (3 * 4)$ na árvore, resultando em 14, e não 20.

Declarações e Tipos

A linguagem suporta a declaração de variáveis com ou sem inicialização imediata. A gramática trata esses dois casos de forma distinta para maior clareza na árvore gerada.

```

1 declaracao:
2     tipo T_ID
3     | tipo T_ID T_ATRIBUICAO expressao
4

```

Ao separar as regras, o analisador consegue diferenciar semanticamente uma simples reserva de memória (apenas declaração) de uma atribuição inicial, gerando nós específicos (`declaracao` ou `declaracao_init`) que facilitam o entendimento do código.

Expressões Matemáticas e Precedência

Para garantir que operações como multiplicação ocorram antes da soma (precedência matemática), a gramática das expressões foi estratificada em três níveis hierárquicos:

- **Fator (Nível mais alto):** Reconhece as unidades básicas, como números, variáveis e expressões entre parênteses.
- **Termo (Nível intermediário):** Processa operações de multiplicação (*) e divisão (/), que têm prioridade sobre a soma.
- **Expressão Simples (Nível base):** Processa somas (+) e subtrações (-).

Essa divisão obriga o analisador a resolver primeiro o que está "mais fundo" na hierarquia (fatores e termos) antes de resolver as somas. Isso elimina a ambiguidade sem a necessidade de código complexo em C, utilizando a própria natureza da gramática para definir a ordem de avaliação.

Exemplos de Análise Sintática

A seguir, são apresentados exemplos de código fonte na linguagem desenvolvida e a respectiva Árvore Sintática gerada pelo analisador sintático.

Inicialização de Variáveis

Na linguagem idealizada, a inicialização de variáveis segue uma lógica similar àquela encontrada em linguagens de programação convencionais. Sob essa perspectiva, os tipos de dados disponíveis são definidos da seguinte maneira:

Inteiro

Equivalente a **int**. É utilizado para representar valores numéricos inteiros.

```

1 // Entrada:
2 inteiro a = 10;
3
4 // Saída:
5 programa
6   lista_comandos
7     declaracao_init
8       tipo
9         inteiro
10        id
11        a
12        literal_inteiro
13        10

```

Quebrado

Equivalente a **float**. É utilizado para representar valores numéricos com casas decimais.

```

1 // Entrada:
2 quebrado b = 20.5;
3
4 // Saída:
5 programa
6   lista_comandos
7     declaracao_init
8       tipo
9         quebrado
10      id
11        b
12      literal_quebrado
13        20.5

```

Texto

Equivalente a **string**. É utilizado para representar sequências de caracteres.

```

1 // Entrada:
2 texto c = "ola mundo";
3
4 // Saída:
5 programa
6   lista_comandos
7     declaracao_init
8       tipo
9         texto
10      id
11        c
12      literal_texto
13        "ola mundo"

```

Caractere

Equivalente a **char**. É utilizado para representar valores de caractere único.

```

1 // Entrada:
2 caractere d = "d";
3
4 // Saída:
5 programa
6   lista_comandos
7     declaracao_init
8       tipo
9         caractere
10      id

```

```

11      d
12      literal_texto
13      "d"

```

Fato

Equivalente a **bool**. É utilizado para representar valores lógicos. Na linguagem, pode assumir os valores **real** e **fake** que representam, respectivamente, **true** e **false**.

```

1 // Entrada:
2 fato x = real;
3 fato y = fake;
4
5 // Saída:
6 programa
7   lista_comandos
8     declaracao_init
9       tipo
10      fato
11      id
12      x
13      literal_fato
14      real
15      declaracao_init
16      tipo
17      fato
18      id
19      y
20      literal_fato
21      fake

```

Serpente

Equivalente a **vetor**. É utilizado para armazenar sequências de elementos do mesmo tipo;

```

1 // Entrada:
2 serpente lista_fatos (fato) = (fake, real, fake);
3
4 // Saída:
5 programa
6   lista_comandos
7     serpente_declaracao
8     lista_fatos
9       tipo
10      fato
11      elementos
12      literal_fato

```

```

13      fake
14      literal_fato
15      real
16      literal_fato
17      fake

```

Estruturas de Controle e Repetição

Condisional Se-Senao

Equivalente a **if-else**. O comando **se** é utilizado para avaliar uma condição e executa o bloco correspondente caso a condição seja verdadeira. O comando **senao** define o bloco alternativo a ser executado quando a condição do **se** não é satisfeita.

```

1 // Entrada:
2 se (a > b) {
3     b = b + 1;
4 }
5 senao {
6     a = a + 1;
7 }
8
9 // Saida:
10 programa
11     lista_comandos
12         se_senao_statement
13             condicao
14                 op_relacional
15                     >
16                     id
17                     a
18                     id
19                     b
20             lista_comandos
21                 atribuicao
22                     id
23                     b
24                 op_aritmetica
25                     +
26                     id
27                     b
28                     literal_inteiro
29                     1
30         bloco_senao
31             lista_comandos
32                 atribuicao
33                     id
34                     a
35                 op_aritmetica

```

```

36      +
37      id
38      a
39      literal_inteiro
40      1

```

Enquanto

Equivalente a **while**. O comando **enquanto** implementa laços de repetição, permitindo que um bloco de código seja executado repetidamente enquanto uma condição lógica for verdadeira.

```

1 // Entrada:
2 enquanto (c >= 15) {
3     c = c - 1;
4 }
5
6 // Saída:
7 \begin{lstlisting}[style=MyLang]
8 programa
9     lista_comandos
10    enquanto_loop
11        condicao
12            op_relacional
13                >=
14                id
15                c
16                literal_inteiro
17                15
18        lista_comandos
19        atribuicao
20            id
21            c
22            op_aritmetica
23                -
24            id
25            c
26            literal_inteiro
27            1

```

Até

Equivalente a **for**. O comando **até** permite definir um contador e executar um bloco de código repetidas vezes até que determinada condição seja satisfeita.

```

1 // Entrada:

```

```

2 ate (i > d) {
3   i = i + 1;
4 }
5
6 // Saída:
7 programa
8   lista_comandos
9     ate_loop
10    condicao
11      op_relacional
12        >
13        id
14        i
15        id
16        d
17   lista_comandos
18     atribuicao
19       id
20       i
21     op_aritmetica
22       +
23       id
24       i
25       literal_inteiro
26       1

```

Entrada e Saída de Dados

Receba

Equivalente a **input**. É utilizado para realizar a entrada de dados, permitindo que o programa capture informações fornecidas pelo usuário durante a execução.

```

1 // Entrada:
2 inteiro x;
3 receba(x);
4
5 // Saída:
6 programa
7   lista_comandos
8     declaracao
9     tipo
10    inteiro
11    id
12    x
13    receba
14    x

```

Imprima

Equivalentе a **print**. É responsável pela saída de dados, permitindo exibir informações na tela durante a execução do programa.

```
1 // Entrada:  
2 inteiro y = 10;  
3 imprima(y);  
4  
5 // Saída:  
6 programa  
7     lista_comandos  
8         declaracao_init  
9             tipo  
10            inteiro  
11            id  
12            y  
13            literal_inteiro  
14            10  
15            imprima  
16            id  
17            y
```

Funções

Na linguagem, as funções podem ser dos respectivos tipos: inteiro, quebrado, nada.

Devolva

Equivalentе a **return**. É um comando de controle que interrompe a execução de uma função e retorna um valor (ou nenhum) ao ponto onde ela foi chamada.

```
1 // Entrada:  
2 devolva a;  
3  
4 // Saída:  
5 programa  
6     lista_comandos  
7         devolva  
8             id  
9             a
```

Funções do tipo Inteiro

São as funções que devolvem um valor do tipo inteiro.

```
1 // Entrada:  
2 inteiro exemplo(inteiro a) {  
3     devolva a;  
4 }  
5  
6 // Saida:  
7 programa  
8     lista_comandos  
9         funcao_definicao  
10        exemplo  
11        tipo  
12        inteiro  
13        parametros  
14        parametro  
15        a  
16        tipo  
17        inteiro  
18        lista_comandos  
19        devolva  
20        id  
21        a
```

Funções do tipo Quebrado

São as funções que devolvem um valor do tipo quebrado.

```
1 // Entrada:  
2 quebrado exemplo(quebrado b) {  
3     devolva b;  
4 }  
5  
6 // Saida:  
7 programa  
8     lista_comandos  
9         funcao_definicao  
10        exemplo  
11        tipo  
12        quebrado  
13        parametros  
14        parametro  
15        b  
16        tipo  
17        quebrado  
18        lista_comandos  
19        devolva  
20        id
```

Funções do tipo Nada

São as funções que não devolvem nada.

```

1 // Entrada:
2 nada exemplo() {
3 }
4
5 // Saída:
6 programa
7 lista_comandos
8     funcao_definicao
9     exemplo
10    tipo
11    nada
12    parametros
13    nenhum
14    bloco_vazio

```

Alterações em Relação ao Trabalho 1

Para viabilizar a análise sintática, foram feitas alterações na camada léxica (Flex) definida no Trabalho 1 :

1. **Remoção de Impressões Diretas:** O Flex deixou de imprimir os tokens (ex: ‘<inteiro>’) diretamente na saída padrão.
2. **Retorno de Tokens:** O Flex agora retorna constantes inteiras (definidas pelo Bison) que representam o tipo do token.
3. **Passagem de Valores:** Foi implementada a lógica de ‘strdup(yytext)’ para preencher a estrutura ‘yylval’, permitindo que o Bison tenha acesso ao texto original dos identificadores e literais para a construção da árvore.

Considerações Finais

O desenvolvimento do analisador sintático com Bison permitiu estruturar formalmente a gramática da linguagem. A abordagem de construção da árvore baseada em manipulação de strings provou-se eficaz para atender aos requisitos de visualização hierárquica solicitados, mantendo o código modular e legível. A gramática livre de contexto elaborada cobre todas as especificações do projeto, incluindo expressões complexas, funções e estruturas de controle aninhadas.

GitHub

Link para acesso ao repositório do projeto: <https://github.com/arthurscorreia/syntax-analysis>.