

UNIVERSIDADE FEDERAL DE VIÇOSA
INSTITUTO DE CIÊNCIAS EXATAS E TECNOLÓGICAS
BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

CCF441 – COMPILADORES
LINGUAGEM TAO – ENTREGA 1 DO TRABALHO PRÁTICO

GRUPO:

3873 – Germano Barcelos dos Santos

3051 – Henrique de Souza Santana

3890 – Otávio Santos Gomes

3877 – Pedro Cardoso de Carvalho Mundim

3495 – Vinícius Júlio Martins Barbosa

Relatório de trabalho prático da disciplina
CCF441 – COMPILADORES, apresentado à
Universidade Federal de Viçosa.

FLORESTAL
MINAS GERAIS - BRASIL
JUNHO 2022

SUMÁRIO

1	Introdução	3
2	Especificação	3
2.1	Nome da Linguagem e Origem	3
2.2	Paradigma de Programação	4
2.3	Léxico	4
2.4	Tipos de Dados	6
2.5	Literais	7
2.6	Comandos	7
2.7	Expressões	12
2.8	Sintaxe	13
3	Análise Léxica	17
4	Alterações	21
5	Yacc	21
6	Considerações Finais	21
	Referências	22

1 INTRODUÇÃO

Nesta primeira parte do trabalho prático da disciplina CCF441 - Compiladores, o objetivo é definir uma linguagem de programação própria, não só prática, mas também criativa. Além disso, também é objetivo do trabalho de implementar o seu analisador léxico, utilizando para isso o Lex e a linguagem de programação C. Dessa forma, espera-se obter um conhecimento mais aprofundado da estrutura de um analisador léxico, além de destacar características de linguagens de programação que geralmente só são observadas por projetistas de uma linguagem. Assim, apresentamos a especificação de uma linguagem, bem como seu analisador léxico a seguir.

2 ESPECIFICAÇÃO

Nesta seção será mostrada a especificação da linguagem, apresentando o nome, origem do nome, tipos de dados, comandos, palavras-chave, dentre outros.

2.1 NOME DA LINGUAGEM E ORIGEM

O nome da linguagem é **Tao**, seu logotipo é mostrado na Figura 1, e tem origem na tradição filosófica e religiosa chinesa do taoismo, que enfatiza a harmonia com o *tao* (, "caminho") (1). Devido a características da linguagem que serão descritas nas seções seguintes, a linguagem Tao busca deixar o programador como responsável por melhor "harmonizar" o código e explicitamente gerenciar tipos e memória.

Três conceitos do taoismo foram incorporados no léxico da linguagem. O primeiro deles é o do *yin* e *yang*. Ambos são forças opostas e complementares, que se originam do *wuji* – um "vazio primordial". O segundo conceito é o dos trigramas. Yin e yang se combinam nos chamados oito trigramas, que representam princípios fundamentais da realidade correlacionados. Como mostra a Figura 2, yin é representado por uma barra quebrada, e yang por uma barra completa. Se considerarmos o yin como valor 0, e o yang como valor 1, e se lermos os trigramas de baixo para cima, podemos associar a cada um deles um valor codificado em binário.



Figura 1 – Logotipo da linguagem.

Além disso, o terceiro conceito, emprestado do *I Ching* (, geralmente traduzido como "O Livro das Mudanças") (2), é que esses oito trigramas também se recombinaem, formando 64 hexagramas. Os hexagramas igualmente podem ser associados a valores binários, apesar dessa associação não ser usada nos textos clássicos. Cada hexagrama também representa uma ideia, um conceito, e buscamos associá-los a construções de programação.

As palavras *wuji*, *yin* e *yang* se tornaram palavras reservadas, e cada trígama e hexagrama também, escritos como sequências de, respectivamente, 3 ou 6 símbolos : para yin ou | para yang. Por exemplo, o trígama do trovão se tornou | : : e o do vento se tornou : | |.

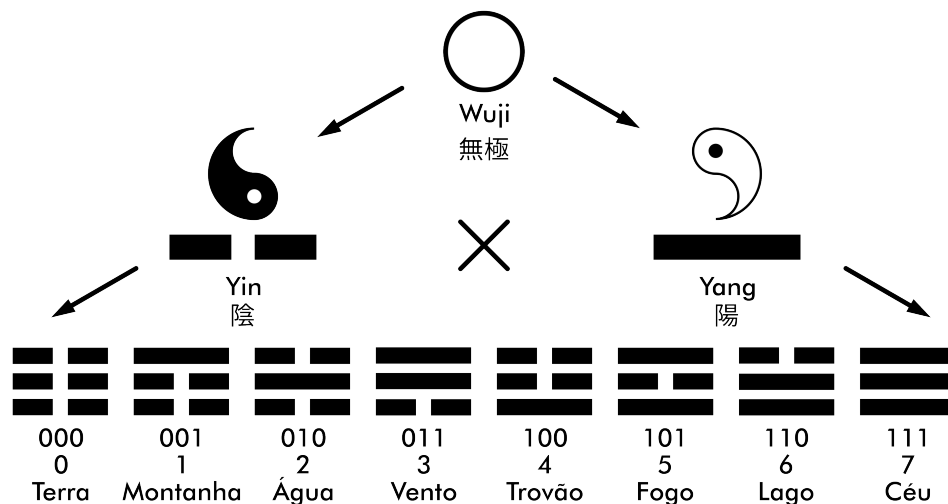


Figura 2 – Os trigramas e sua origem.

2.2 PARADIGMA DE PROGRAMAÇÃO

Como requisito na especificação do trabalho, o paradigma de programação da linguagem é **procedural**. Além disso, tanto quanto possível, buscamos tornar a linguagem orientada a expressões, de forma que, por exemplo, atribuição e condicional (*if*) podem ser usados como expressões. Deixamos em aberto a possibilidade de definir funções aninhadas e passá-las como parâmetro para outras funções, incorporando aspectos do paradigma funcional. Por fim, pretendemos implementar mecanismos simples de polimorfismo de coerção, de sobrecarga e paramétrico.

2.3 LÉXICO

Um primeiro aspecto importante do léxico da linguagem se refere aos identificadores, que podem ser comuns, próprios ou simbólicos. Identificadores próprios devem começar com uma letra maiúscula, seguida de uma quantidade indefinida de letras, dígitos ou traços baixos (_). Os comuns seguem a mesma regra, mas começam com letra minúscula ou traço baixo. Identificadores simbólicos contêm uma ou mais repetições dos caracteres !\$%&*+. /<=>? | :#@^~-. .

Todo identificador, seja comum, próprio ou de símbolo, pode ser prefixado com um escopo, caso seja importado de algum módulo. Esse prefixo é uma sequência de um ou mais identificadores próprios seguidos de um ponto final. Por exemplo, uma função chamada `function` importada do módulo `Some.Module` pode ser referida como `Some.Module.function`. Identificadores escritos assim são chamados identificadores qualificados, e são considerados um único lexema. Identificadores comuns nomeiam variáveis, funções e procedimentos, identificadores simbólicos nomeiam operadores, e identificadores próprios nomeiam tipos e módulos.

Por simplicidade e também clareza de código desenvolvido na linguagem, optamos por definir todas as palavras-chave como palavras reservadas, evitando assim ambiguidades. Na verdade, como supracitado, todos os hexagramas e trigramas são palavras reservadas, mas nem todos os hexagramas tem um significado atribuído (ainda). Dessa forma, é possível apontar as seguintes palavras-chave:

- wuji – o wuji representa o vazio primordial, então foi escolhido para definir procedimentos, o análogo em C para função de retorno *void*.
- yin – o yin é a força negativa/passiva, então escolhemos esta palavra-chave para criar as entidades de programação que armazenam valores, ou seja, definir variáveis.
- yang – o yang é a força positiva/ativa, então escolhemos esta palavra-chave para criar as entidades de programação que produzem valores, ou seja, definir funções.

Em seguida, temos os trigramas, que são associados a conceitos mais específicos. Junto aos trigramas temos os hexagramas, que são agrupados de acordo com seu prefixo. Por exemplo, o hexagrama `:|:|:` possui o prefixo do trigrama da água (`:|:`), então será usado em construções relacionadas a este trigrama. Nas Seções 2.6 e 2.7 serão mostrados os hexagramas usados em cada construção específica.

- `:::` – representa receptividade, está relacionado a definição e declaração de tipos;
- `::|` – representa estabilidade, está relacionado a alocação de memória;
- `:|:` – representa movimento, está relacionado a condicionais;
- `:||` – representa pervasividade, está relacionado a repetição;
- `|::` – representa perturbação, está relacionado a liberação de memória;
- `|:|` – representa iluminação, está relacionado a verificação de tipos e valores;
- `||:` – representa abertura, está relacionado a importação de módulos;
- `|||` – representa criação, está relacionado a exportação de módulos.

Além dessas palavras reservadas, alguns símbolos são reservados para controlar a sintaxe da linguagem. Isso é importante, pois, como veremos na Seção 2.7, o programador pode definir novos operadores. Embora `@`, `:`, `=` e `.` possam ser usados como identificadores de operadores – desde que acompanhados de outros símbolos –, os demais símbolos da lista a seguir *nunca* podem ser usados como operadores.

- Todo comando deve encerrar com ponto e vírgula (`;`) ou com nova linha, então esses lexemas representam o mesmo token de fim de linha;
- Caracteres arroba (`@`) prefixam nomes de tipo, indicando que trata-se de ponteiro;
- Dois pontos (`:`) indicam o tipo da variável que o precede;
- Sinal de igual (`=`) é usado em atribuições e definições;

- Ponto final (.) é usado para acessar campos em estruturas e módulos;
- Vírgula (,) é usada como separador em listas de parâmetro;
- Colchetes ([]) são usados para acessar posições em ponteiros;
- Parênteses (()) aninham expressões, e chaves ({}) aninham blocos.
- Comentários de uma única linha são indicados com dois apóstrofes ('), e comentários multilinha são indicados com três apóstrofes no começo e três apóstrofes no final (''').

2.4 TIPOS DE DADOS

A linguagem Tao possui os seguintes tipos primitivos:

- Int – valores inteiros de 32 bits com sinal;
- Long – valores inteiros de 64 bits com sinal;
- Real – valores de ponto flutuante de 64 bits;
- Char – valores de caractere de 8 bits;
- Bool – valores booleanos;
- @Type – ponteiro para Type, sendo Type qualquer outro tipo.

Uma observação importante é quanto a notação de tipos para funções, não inclusa na lista acima. Quando uma variável na verdade guarda um ponteiro para função, no lugar da anotação de tipo normalmente usada, deve ser usado o hexagrama `|||::|` seguido de uma lista de tipos, e opcionalmente seguido novamente pelo mesmo hexagrama e um único tipo. Quando o hexagrama aparece apenas uma vez, trata-se de um ponteiro para um subprograma que espera como parâmetro os tipos listados. Quando o hexagrama aparece duas vezes, trata-se de um ponteiro para uma função que espera como parâmetro os tipos listados depois do primeiro hexagrama, e retorna valores do tipo mencionado depois do segundo hexagrama. Por exemplo, `|||::| Int, Int |||::| Bool` denota funções que recebem dois inteiros e retorna booleano, `|||::| @Char` denota um subprograma que recebe um ponteiro para caractere.

Além disso, é possível definir tipos compostos, usando uma combinação dos mecanismos de produto cartesiano e união disjunta. Ao definir novos tipos, o programador pode primeiramente definir uma lista de construtores para aquele tipo (união disjunta), e cada construtor pode ter zero ou mais campos nomeados, de tipos arbitrários (produto cartesiano). O Algoritmo 2.1 mostra exemplos dessa definição.

```
1  ' ' três construtores, nenhum campo
2  ::: Color = Red, Green, Blue
3
4  ' ' um construtor com dois campos
5  ::: Person = Person(name: @Char, age: Int)
6
7  ' ' três construtores, com campos variados
8  ::: Tree = Node(x: Int, lt: @Tree, rt: @Tree), Leaf(x: Int), Empty
```

Algoritmo 2.1 – Tipos compostos.

2.5 LITERAIS

Cada tipo primitivo possui literais associadas. Valores inteiros são expressados por sequências de dígitos decimais, ou por dígitos hexadecimais, prefixados com `0x`. Valores reais são expressados com dígitos decimais separados com um ponto, seguidos opcionalmente por um expoente em notação científica.

Caracteres são expressados por dois apóstrofes com um caractere ou uma sequência de escape dentro. Sequências de escape começam com contrabarra (`\`) e são seguidas ou de um valor hexadecimal arbitrário entre `00` e `ff`, (ex.: `'\x20'`), ou então são sequências bem conhecidas como `'\n'` de nova linha ou `'\t'` de tabulação. Para expressar o próprio apóstrofo também é necessário esse escape (`'\''`). Apesar de não haver strings como tipos primitivos, strings são vetores de caracteres, como em C, que usa ponteiros, e podem ser expressas em literais. As regras são as mesmas para literais de caracteres, porém ao invés de apóstrofes, usa-se aspas, a string pode ser vazia, e ao invés de ser necessário usar escape no apóstrofo, deve ser usado o escape nas aspas.

Literais booleanas não são exatamente literais, mas sim uma expressão de construção de tipo. O tipo `Bool` é definido como `::: Bool = False, True`, e portanto seus valores são usados como `False` e `True`.

2.6 COMANDOS

Os comandos disponíveis são: exportação de módulo, importação, definição de variável, definição de função, definição de operador, definição de procedimento, definição de tipo, declaração de tipo, chamada de procedimento, condicional, casamento de tipos, repetição (*while* e *repeat*) e liberação de memória.

O Algoritmo 2.2 só aparece opcionalmente no topo do arquivo e não se repete. Utilizamos o trigrama `|||` para denotar a exportação e o hexagrama `|||:::` para listar quais procedimentos, funções, variáveis ou tipos serão exportadas.

Com relação à importação de módulos, utilizamos o trigrama `||:` para denotar a importação, o hexagrama `||::||` para especificar o que se deseja importar de dentro de um módulo. Além disso, o hexagrama `||::||:` pode ser usado para renomear determinado módulo

```

1  '' arquivo1.tao
2  ||| This.Module.Name |||::: proc1, func2, var4
3
4  '' arquivo2.tao
5  ||| Another.Module |||::: var8, proc16, Type32

```

Algoritmo 2.2 – Comando de exportação de módulo.

```

1  '' equivalente em Python:
2  '' from src.patricia import *
3  ||: Src.Patricia
4
5  '' from math import sin, cos, tan
6  ||: Math ||::|| sin, cos, tan
7
8  '' import sqlite3 as sql
9  ||: SQLite3 ||::||: SQL

```

Algoritmo 2.3 – Comando de importação.

importado, facilitando sua utilização posteriormente no código (Alg. 2.3). Nesse último caso, os identificadores pertencentes a esse módulo serão acessados como identificadores qualificados (ex.: `SQL.connect`).

```

1  ::::: String = @Char
2
3  ::: Person = Person(name: @Char, age: Int)
4
5  yin pi: Real = 3.1415926535
6
7  yang initial(p: Person): Char = p.name[0]
8
9  '' ::::|| lista parâmetros de tipo para
10 '' a função/procedimento/operador que vier em seguida
11 ::::|| T
12 yang len(list: @List(T)): Int = { ''' expressão ''' }
13
14 wuji main() { }

```

Algoritmo 2.4 – Comando de declaração de tipo, e comandos de definições: tipos, variáveis, funções e procedimentos.

O Algoritmo 2.4 apresentam os comandos de definição e declaração. A declaração de tipo é dada pelo hexagrama `:::::`, que atribui um novo nome a um tipo já existente. Já com o trigrama `:::`, é possível declarar novos tipos de dados, como especificado na Seção 2.4. `yin` é

utilizado em uma definição de variável. Nessa definição, após o comando `yin` aponta-se o nome da variável seguido de dois pontos e seu tipo, podendo já ser inicializada com um valor.

`yang`, por sua vez, é o comando responsável pela a definição de funções. Utiliza-se o comando `yang` seguido do nome da função, seus parâmetros e tipos dentro de parênteses. Em sequência, coloca-se `:` e o tipo de retorno da função. Por fim, utiliza-se um `=` e o código pode estar em sequência entre chaves, caso seja multilinhas, ou apenas a linha a ser executada logo após o `=`. Por fim, `wuji` pode ser utilizado para definir subprogramas, que são a mesma coisa que funções, porém não retornam nenhum valor e sua chamada não pode ser usada como expressão.

No corpo de subprogramas e funções, é possível usar o comando `||||:`, seguido de uma expressão, para encerrar a execução imediatamente e retornar o valor dessa expressão.

```

1  '' operadores associativos à esquerda com :::| |:
2  '' número indica a precedência
3  yang :::| |: 5 <<(x: Real, exp: Real): Real = { '' expressão '' }
4  yang :::| |: 5 >>(x: Real, exp: Real): Real = { '' expressão '' }
5
6  '' operadores associativos à direita com :::| :|
7  yang :::| :| 5 +>(item: Int, head: @List): @List = {
8      '' expressão ''
9  }
10
11 '' operadores não associativos com :::| ::
12 yang :::| :: 4 ==(s0: @Char, s1: @Char): Bool = {
13     '' expressão ''
14 }
```

Algoritmo 2.5 – Comando de definição de operadores.

Veja no Algoritmo 2.5 que é possível também definir operadores. Essa definição se dá como na definição de uma função. No entanto, é preciso indicar a precedência do operador. Além disso, deve-se ter exatamente 2 parâmetros. Dessa forma, define-se um operador da seguinte forma: utiliza-se o comando `yang` seguido do hexagrama `:::| |:`, `:::| :|` ou `:::| ::`, seguido de sua precedência. Esses hexagramas indicam a associatividade do operador (esquerda, direita ou nenhuma respectivamente). Logo em sequência, apresenta-se o nome do operador, seguido de qual o nome e o tipo dos seus 2 parâmetros. Por fim, apresenta-se seu valor de retorno e enfim descreve-se a sua execução no bloco de código.

No Algoritmo 2.6 mostramos que a chamada de um procedimento é bastante simples, podendo-se passar quaisquer parâmetros, conforme estabelecido na definição do procedimento. A sintaxe de chamada de procedimento e de função é exatamente igual, porém futuramente na análise sintática é que deve ser conferido se uma chamada de procedimento não está sendo usada como expressão.

Veja no Algoritmo 2.7 a implementação do comando condicional. Nele, o trigrama `:| :`

```

1 wuji dfs(node: @Tree) { }
2
3 yin tree: @Tree
4 dfs(tree)

```

Algoritmo 2.6 – Comando de chamada de procedimento.

```

1 |:|: cond |:|:: {
2     '' executa se cond == True
3 }
4
5 |:|: cond |:|:: {
6     '' executa se cond == True
7 } |:|::|| {
8     '' executa se cond == False
9 }
10
11 |:|: cond |:|:: {
12     '' executa se cond == True
13 } |:|::|: cond2 |:|:: {
14     '' executa se cond == False e cond2 == True
15 } |:|::|| {
16     '' executa caso contrário
17 }

```

Algoritmo 2.7 – Comando condicional.

representa o tradicional comando *se*, e o hexagrama $|:|:|:$ aparece depois da condição. *Senão* é representado pelo hexagrama $|:|::||$, e aparece opcionalmente. Uma quantidade indefinida de *senão-ses* pode aparecer depois da primeira condição, indicados pelo hexagrama $|:|::|:$.

```

1 ::: Tree(K) = Node(x: K, lt: @Tree(K), rt: @Tree(K)), Leaf(x: K)
2
3 yin tree: @Tree(Int) = '' expressão
4 |:| @tree
5 |:|||| Node(x, l, r) |:|:|: { ''executa se for Node'' }
6 |:|||| Leaf(x) |:|:|: { ''executa se for Leaf'' }
7 |:|:|: { ''default'' }

```

Algoritmo 2.8 – Comando de casamento de tipo.

No Algoritmo 2.8 mostramos a estrutura de um comando similar a um `switch-case` da linguagem C. O comando é definido pelo trigramma $|:|$ seguido de qual expressão será utilizada na análise feita pelo comando. O hexagrama $|:||||$ serve para realizar o casamento, caso seja verdadeiro, será executado o código presente em sequência, depois do hexagrama $|:|:|:$.

Por fim, o hexagrama `:|:|:|:` representa o caso padrão, que será executado quando nenhuma possibilidade anterior for executada.

```

1  :| cond :|:|:| {
2      '' executa enquanto cond == True
3      :|:|:|: '' pula para a próxima iteração, comando "continue"
4  }
5
6  :| cond :|:|:|: { '' executa depois do laço'' } :|:|:| {
7      '' executa enquanto cond == True
8  }
```

Algoritmo 2.9 – Comando de repetição com teste antes do laço.

No Algoritmo 2.9, temos os comandos responsáveis pela repetição. O trigramma `:|` representa o início do comando de repetição. Ele é sempre seguido de uma condição. Depois pode ser seguido de diferentes hexagramas. Se seguido de `:|:|:|`, executará enquanto o comando descrito for verdadeiro. Se, por sua vez, a condição for seguida do hexagrama `:|:|:|:`, este será seguido de um bloco de código a ser executado após a finalização do laço. Por fim, deve ser seguido do hexagrama `:|:|:|` que terá um trecho de código que será executado enquanto a condição apresentada for verdadeira.

```

1  :|:|:| {
2      '' executa até que cond == True
3      :|:|:|: '' sai imediatamente do laço, comando "break"
4  } :|:|:| cond
5
6  :|:|:| {
7      '' executa até que cond == True
8  } :|:|:| cond :|:|:|: { '' executa depois do laço'' }
```

Algoritmo 2.10 – Comando de repetição com teste depois do laço.

No Algoritmo 2.10, por sua vez, utilizamos o hexagrama `:|:|:|` seguido de um bloco de código, seguido do hexagrama `:|:|:|` e da condição de execução desse trecho de código. Esse código será executado até que a condição se torne verdadeira, porém ao contrário da estrutura apresentada anteriormente, executará pelo menos uma vez, sendo checada a condição apenas após a execução.

Em ambos os comandos, é possível usar os hexagramas `:|:|:|:` e `:|:|:|` para controlar o fluxo da repetição, sendo que o primeiro é o comando que sai do laço imediatamente (break), e o segundo é o que pula para uma próxima iteração do laço (continue).

No Algoritmo 2.11, observe que utilizamos o trigramma `|:|:` para a liberação de memória. Fazendo um paralelo com a linguagem C, funcionaria como uma chamada à função `free()`.

```

1 | yin list: @List = ' ' expressão
2 | :: list

```

Algoritmo 2.11 – Comando de liberação de memória.

Precedência	Operadores	Associatividade
8	**	Direita
7	* / & ^	Esquerda
6	+ -	Esquerda
5	<< >>	Esquerda
4	== != < <= >= >	
3	&&	Esquerda
2		Esquerda
1	(nenhum)	

Tabela 1 – Operadores predefinidos.

2.7 EXPRESSÕES

Toda expressão pode ser aberta numa sequência de expressões e comandos, usando chaves ({ }), e, nesse caso, o resultado da última expressão do bloco se torna o resultado do todo. Expressões também podem ser aninhadas com parênteses.

A proposta da linguagem permite ao programador definir seus próprios operadores, mas já possui operadores predefinidos (que podem inclusive ser sobrecarregados), além de algumas expressões especificadas com os trigramas e hexagramas. A Tabela 1 mostra os operadores predefinidos, suas precedências e suas associatividades. Quanto maior o valor de precedência, maior a prioridade. Os operadores matemáticos +, -, *, /, ==, !=, <, <=, >= e > têm seu significado intuitivo. O operador ** se refere à exponenciação. Os operadores &, |, ^, << e >> são binários, e se referem respectivamente às operações *and*, *or*, *xor*, *shift-left* e *shift-right*. Por fim, os operadores lógicos && e || também têm seu significado intuitivo.

No total, as expressões disponíveis são: condicional, casamento de tipo, atribuição (=), acesso a campo (.comid), acesso a posição ([i]), acesso direto (@), referência (\$), negação binária (~), negação lógica (!), negação numérica (-), operações infixadas com ordem de precedência de 1 a 8, literais, variáveis, alocação de memória, construção de tipo composto e chamada de função.

As expressões de casamento de tipo e a condicional seguem a mesma sintaxe de quando são usadas como comandos, com a diferença de que o "else" do condicional é obrigatório. As expressões de negação e de manipulação de ponteiro são os únicos operadores prefixados. Exemplos de cada expressão são mostrados no Algoritmo 2.12.

```

1  '' expressão de atribuição, e diferentes níveis de operações infixadas
2  x = y = 3 | 2 ** 3 >> 0xA + 0xf - 23
3  '' equivalente à seguinte expressão:
4  x = (y = ((3 | (2 ** 3)) >> ((0xA + 0xf) - 23)))
5
6  '' negações
7  cond = !(-x < ~y)
8
9  '' manipulação de ponteiros
10 y = @($x + 10)
11
12 '' acesso a posição, seguido de acesso a campo
13 people[10].name = "Alfred Aho"
14
15 '' alocação de memória no monte
16 yin p: @Person = ::| Person
17 '' alocação de memória na pilha
18 p = ::||:: Person
19 '' alocação de 10 posições no monte
20 p = ::| Person ::|:| 10
21 '' alocação de 10 posições na pilha
22 p = ::||:: Person ::|:| 10
23 '' construção de tipo
24 @p = Person("Fulano", 42)
25 '' alocação + construção
26 p = ::| Person("Fulano", 42)
27 '' alocação + construção de 20 posições
28 p = ::| Person("", 0) ::|:| 20
29
30 '' chamada de função
31 x = a + random() * (b - a)

```

Algoritmo 2.12 – Expressões.

2.8 SINTAXE

A gramática da linguagem será expressada na Forma de Backus-Naur (BNF) (3). Nessa notação, as variáveis da gramática são expressas <desta-forma>, as produções são indicadas por $::=$, e terminais são indicados entre aspas. Evidentemente, todo terminal seria representado por um token, mas, para facilitar a representação de tokens cujos lexemas podem variar, estendemos essa notação usando parênteses, por exemplo (LiteralInt) se refere uma literal de inteiro.

Especificamente, esses são os tokens para os quais usamos essa notação:

- (LiteralInt) – literal de número inteiro;
- (LiteralReal) – literal de número real;

- (LiteralChar) – literal de caractere;
- (LiteralString) – literal de string;
- (Endl) – fim de linha;
- (TrigN) – trigrama de valor N, ex.: (Trig6) é ||::;
- (HexN) – hexagrama de valor N, ex.: (Hex42) é |:|:|:|::;
- (ComID) – identificador comum;
- (ProID) – identificador próprio;
- (SymID) – identificador simbólico;
- (QComID) – identificador comum qualificado;
- (QProID) – identificador próprio qualificado;
- (QSymID) – identificador simbólico qualificado.

Assim, segue a GLC da linguagem, segmentada em partes relacionadas para facilitar seu entendimento. Como convencional, a primeira variável mostrada (<program>) é a variável de partida.

```

1 <program> ::= <module-decl> <top-stmts>
2
3 <module-decl> ::= {Trig7} <pro-id> {Hex56} <exports> {Endl} | ""
4 <exports> ::= <export> ", " <export-id> | <export-id>
5 <export-id> ::= {ComID} | {ProID} | {SymID}
6
7 <top-stmts> ::= <top-stmts> {Endl} <top-stmt> | <top-stmt>
8 <top-stmt> ::= <import>
9               | <callable-def>
10              | <def-type-params> <callable-def>
11              | <var-def>
12              | <type-def>
13              | <type-alias>
14              | ""

```

Algoritmo 2.13 – Início da GLC.

No Algoritmo 2.13 são mostradas as construções que devem aparecer no topo de um programa, antes de entrar em qualquer bloco. Como supracitado, a declaração de módulo é opcional e necessariamente deve ser a primeira linha do programa, representada pela variável <module-decl>. A variável <top-stmts> é usada para permitir que apenas comandos de definição e declaração possam ser usados fora de algum bloco.

Em seguida, a variável `<stmts>` (Alg. 2.14) é a que permite quaisquer comandos, e aparecerá no corpo de outras produções mais a frente.

```

1  <stmts> ::= <stmts> (Endl) <stmt> | <stmt>
2  <stmt> ::= <import>
3           | <callable-def>
4           | <def-type-params> <callable-def>
5           | <var-def>
6           | <type-def>
7           | <type-alias>
8           | <call>
9           | <if-stmt>
10          | <match>
11          | <while>
12          | <repeat>
13          | <free>
14          | <break>
15          | <continue>
16          | <return>
17          | <expr>
18          | ""

```

Algoritmo 2.14 – Comandos expressos na GLC.

O Algoritmo 2.15 mostra as variáveis que expressam as construções de definição e declaração de tipos, variáveis, funções, etc., além do comando de importação.

No Algoritmo 2.16 temos alguns comandos e expressões que podem ser usadas como comandos. No caso, o comando condicional, com o detalhe que sua sintaxe é ligeiramente diferente quando usado como expressão. Isso pode ser visto no corpo das produções referentes às variáveis `<if-stmt>` e `<if-expr>`. O comando de casamento de tipos não muda a sintaxe quando usado como expressão, então não foi necessário distinguir. Comandos auxiliares de controle de fluxo também aparecem aqui, com os hexagramas que correspondem ao `break`, `continue` e `return`.

As expressões na GLC requerem uma atenção maior (Alg. 2.17), devido aos diferentes níveis de precedência e associatividade. Depois de permitir a abertura de expressões em blocos, depois da possibilidade de usar expressões condicionais, de casamento de tipos, de atribuição, de endereçamento e operações unárias prefixadas, vêm os operadores infixados. As variáveis `<expr-N>`, `<expr-n-N>`, `<expr-l-N>` e `<expr-r-N>` na verdade representam as variáveis `<expr-1>`, `<expr-n-1>`, `<expr-l-1>`, `<expr-r-1>` até `<expr-8>`, `<expr-n-8>`, `<expr-l-8>` e `<expr-r-8>`, e onde aparece `N+1` no nome das variáveis e tokens no corpo de suas produções, deve-se trocar pelo valor de `N + 1`. Por exemplo, a variável `<expr-l-8>` produz `<expr-l-8>` (SymID18) `<expr-9>`.

Primeiramente, a diferenciação entre `<expr-N>`, `<expr-n-N>` e `<expr-l-N>` serve para

```

1  <import> ::= (Trig6) <pro-id>
2             | (Trig6) <pro-id> (Hex51) <export-list>
3             | (Trig6) <pro-id> (Hex54) (ProID)
4
5  <def-type-params> ::= (Hex03) <type-params>
6  <callable-def> ::= <func-def> | <op-def> | <proc-def>
7
8  <var-def> ::= "yin" (ComID) ":" <type-id>
9             | "yin" (ComID) ":" <type-id> "=" <expr>
10
11 <func-def> ::= "yang" (ComID) "(" <param-list> ")" ":" <type-id> "=" <expr>
12
13 <op-def> ::= "yang" <op-assoc> <op-prec> (SymID) "(" <param> "," <param>
14           ↪ ")" ":" <type-id> "=" <expr>
15 <op-assoc> ::= (Hex04) | (Hex05) | (Hex06)
16 <op-prec> ::= "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8"
17
18 <proc-def> ::= "wuji" (ComID) "(" <param-list> ")" <block>
19
20 <type-def> ::= (Trig0) (ProID) <type-param-list> "=" <constructors>
21 <constructors> ::= <constructors> "," <constructor> | <constructor>
22 <constructor> ::= (ProID) "(" <param-list> ")"
23
24 <type-alias> ::= (Hex00) (ProID) <type-param-list> "=" <type-id>
25
26 <param-list> ::= <params> | ""
27 <params> ::= <params> "," <param> | <param>
28 <param> ::= (ComID) ":" <type-id>

```

Algoritmo 2.15 – Comandos de definições e declarações na GLC.

indicar a associatividade, fazendo com que a árvore de derivação da gramática siga a associatividade do operador, e também impedindo que associatividades diferentes se misturem num mesmo nível de precedência. Outro detalhe é que os tokens (SymIDxN) não aparecem na legenda mostrada no começo desta seção. Isso é porque, pensando numa próxima entrega do trabalho, será necessário que o analisador sintático seja capaz de consultar a tabela de símbolos e verificar, para um dado operador (token (SymID)), se possui a associatividade e precedência esperadas.

Chegando na variável <expr-9> se encerra os níveis de precedência e enfim a expressão pode se tornar uma literal, um identificador, uma alocação de memória, uma construção de tipo uma chamada de função, ou um aninhamento de expressão com parênteses, recomeçando todo o ciclo. No Algoritmo 2.18 vemos as produções referentes a essas últimas construções.

Ao final da gramática (Alg. 2.19), há algumas variáveis que foram usadas em várias produções não necessariamente relacionadas, com destaque à variável <type-id>, que define como um tipo deve ser descrito, e foi uma variável bastante utilizada.


```

1 <if-expr> ::= (Trig2) <expr> (Hex20) <expr> <elif> <else>
2 <if-stmt> ::= (Trig2) <expr> (Hex20) <expr> <elif> <else-stmt>
3 <else-stmt> ::= <else> | ""
4 <else> ::= (Hex19) <expr>
5 <elif> ::= <elif> (Hex18) <expr> (Hex20) <expr> | ""
6
7 <match> ::= (Trig5) <expr> <cases> <default>
8 <cases> ::= <cases> <case> | <case>
9 <case> ::= (Hex47) <case-cond> (Hex42) <expr>
10 <default> ::= (Hex44) <expr> | ""
11 <case-cond> ::= <literal> | <decons>
12 <decons> ::= <pro-id> "(" <com-id-list> ")"
13 <com-id-list> ::= <com-ids> | ""
14 <com-ids> ::= <com-ids> "," (ComID) | (ComID)
15
16 <while> ::= (Trig3) <expr> <step> (Hex31) <block>
17 <repeat> ::= (Hex27) <block> (Hex25) <expr> <step>
18 <step> ::= (Hex28) <block> | ""
19
20 <free> ::= (Trig4) <addr>
21
22 <break> ::= (Hex30)
23 <continue> ::= (Hex26)
24 <return> ::= (Hex62) <expr>

```

Algoritmo 2.16 – Comandos condicional, repetições e liberação de memória na GLC.

3 ANÁLISE LÉXICA

Nesta seção, será mostrado a implementação e desenvolvimento do analisador léxico da linguagem em detalhes. Como dito na introdução, a ferramenta utilizada foi o Lex (4), juntamente com a linguagem C. O arquivo-fonte para o Lex está anexado nos arquivos enviados dentro da pasta /src. Abaixo estão exemplos de entradas e comentários sobre suas respectivas saídas.

As três primeiras linhas do código do Algoritmo 3.2 são os comandos para compilar e executar o analisador lendo o arquivo `match.tao`, e em seguida temos a sua saída. Alternativamente, pode ser executado o comando `make && make run INPUT=match.tao`, usando o arquivo `makefile` preparado.

Primeiramente, é importante observar que organizamos o analisador para fazer a contagem das linhas em que o token analisado está. Com isso, podemos observar que a linguagem ignora o texto que é escrito entre apóstrofes triplos, sendo eles utilizados como comentários, por isso vemos que a linha 1 foi pulada, e, em seguida, foi identificado na linha 2 uma quebra de linha, representada por um (`Endl`). Em seguida, na linha 3, vemos identificações de termos

```

1  <expr> ::= "{" <stmts> (Endl) <expr> "}"
2          | <if-expr>
3          | <match>
4          | <assign>
5          | <expr-addr>
6  <expr-addr> ::= <expr-addr> "[" <expr> "]"
7                | <expr-addr> "." (ComID)
8                | <expr-unary>
9  <expr-unary> ::= <unary-ops> <expr-1> | <expr-1>
10 <unary-ops> ::= <unary-ops> <unary-op>
11 <unary-op> ::= "@" | "$" | "~" | "!" | "-"
12
13 <expr-N> ::= <expr-l-N> | <expr-r-N> | <expr-n-N> | <expr-N+1>
14 <expr-n-N> ::= <expr-N+1> (SymIDnN) <expr-N+1>
15 <expr-l-N> ::= <expr-l-N> (SymIDlN) <expr-N+1>
16 <expr-r-N> ::= <expr-N+1> (SymIDrN) <expr-r-N>
17
18 <expr-9> ::= <literal>
19           | <com-id>
20           | <malloc>
21           | <build>
22           | <call>
23           | "(" <expr> ")"

```

Algoritmo 2.17 – Expressões na GLC.

```

1  <assign> ::= <addr> "=" <expr>
2  <addr> ::= <addr> "[" <expr> "]"
3           | <addr> "." (ComID)
4           | <addr-0>
5  <addr-0> ::= <pointers> <addr-1>
6  <addr-1> ::= <com-id> | "(" <addr> ")"
7
8  <malloc> ::= <malloc-type> <type-id> <malloc-n>
9           | <malloc-type> <expr>
10 <malloc-type> ::= (Trig1) | (Hex12)
11 <malloc-n> ::= (Hex11) (LiteralInt) | ""
12
13 <build> ::= <pro-id> | <pro-id> "(" <exprs> ")"
14 <call> ::= <com-id> "(" <expr-list> ")"
15 <expr-list> ::= <exprs> | ""
16 <exprs> ::= <exprs> "," <expr>

```

Algoritmo 2.18 – Expressões na GLC.

```

1 <block> ::= "{" <stmts> "}" | <stmt>
2
3 <pro-id> ::= (ProID) | (QProID)
4 <com-id> ::= (ComID) | (QComID)
5
6 <type-id> ::= <pointers> <pro-id> <type-param-list>
7               | <pointers> (Hex57) <type-param-list> (Hex57) <type-id>
8               | <pointers> (Hex57) <type-param-list>
9 <type-param-list> ::= "(" <type-params> ")" | ""
10 <type-params> ::= <type-params> "," <type-id> | <type-id>
11 <pointers> ::= <pointers> "@" | ""
12 <literal> ::= (LiteralChar) | (LiteralString) | (LiteralInt) |
   ⇨ (LiteralReal)

```

Algoritmo 2.19 – Variáveis auxiliares da GLC.

```

1 '''A linguagem ignora os comentários
2   escritos entre apóstrofes'''
3 yang add(item: T, tree: @Tree(K, T)): @Tree(K, T) = {
4   yin itemKey: K = getKey(item)
5   yin newLeaf: @Tree(K, T) = ::| Tree(K, T)
6   |:| @tree
7   |:|||| Node(k, l, r) |:|:|:
8   ' ' :|: = if
9   ' ' :|:|: = then
10  ' ' :|:|: = elif
11  ' ' :|:|: = else
12  :|: itemKey < k :|:|:
13    tree.lt = add(item, l)
14  :|:|:|
15    tree.rt = add(item, r)
16  |:|||| Leaf(i) |:|:|: {
17    yin leafKey: K = getKey(i)
18    :|: itemKey < leafKey :|:|:
19      ::| Node(leafKey, newLeaf, tree)
20    :|:|:|
21      ::| Node(itemKey, tree, newLeaf)
22  }
23 }

```

Algoritmo 3.1 – Código do arquivo match.tao.

```

1 $ flex lex.l
2 $ gcc lex.yy.c
3 $ ./a.out < match.tao
4 2@(Endl)
5 3@(Yang)
6 3@(ComId, `add`)
7 3@(Special, `(`)
8 3@(ComId, `item`)
9 ...

```

Algoritmo 3.2 – Compilando o arquivo `lex.l` e trecho da saída para `match.tao`.

como (Yang), que define funções, e `add`, que é um identificador comum, (*ComID*).

```

1 ...
2 6@(Trig, 5)
3 6@(SymPointer)
4 6@(ComId, `tree`)
5 6@(Endl)
6 7@(Hex, 47)
7 7@(ProId, `Node`)
8 7@(Special, `(`)
9 7@(ComId, `k`)
10 7@(Special, `,`)
11 7@(ComId, `l`)
12 7@(Special, `,`)
13 7@(ComId, `r`)
14 7@(Special, `)` )
15 7@(Hex, 42)
16 7@(Endl)
17 ...

```

Algoritmo 3.3 – Trecho da saída de `match.tao` sobre as linhas 6 e 7.

Analisando o trecho referente às linhas 6 e 7 do arquivo *match.tao* (Alg. 3.3), é possível notar a identificação do trigrama relacionado à verificação de tipos e valores (`|:|`, (Trig3)), seguido de um símbolo de ponteiro (`@`, (SymPointer)) e um identificador comum. Na linha abaixo é identificado o hexagrama `|:||||`, identificador próprio, símbolos especiais de parênteses e vírgulas, identificadores comum, e por fim outro hexagrama (`|:|:|:`) e quebra de linha.

Observando outro exemplo (Alg. 3.4), podemos ver como é útil a contagem de linha também para apontar erros no código. Primeiramente, podemos destacar que, na linha 1, a literal de inteiro foi devidamente convertida para seu valor numérico. Depois, na linha 4, na tentativa de definir um operador, foi usado um símbolo não previsto no léxico, um acento grave (```).

```

1 | yin x: Int = 0xaff
2 |
3 | '' operador estranho
4 | yang :::|:: `` (x: Real, y: Real): Int = 2/(1/x + 1/y)

```

Algoritmo 3.4 – Código do arquivo error.tao.

```

1 | 1@(Yin)
2 | 1@(ComId, `x`)
3 | 1@(SymTypeTag)
4 | 1@(ProId, `Int`)
5 | 1@(SymAssign)
6 | 1@(LiteralInt, 2815)
7 | 1@(Endl)
8 | 2@(Endl)
9 | 4@(Yang)
10 | 4@(Hex, 4)
11 | Lexical unbalance at line 4: unknown symbol ``
12 | Lexical unbalance at line 4: unknown symbol ``
13 | 4@(Special, `(`)
14 | ...

```

Algoritmo 3.5 – Trecho da saída para error.tao.

4 ALTERAÇÕES

5 YACC

6 CONSIDERAÇÕES FINAIS

Consideramos que a especificação de uma linguagem foi uma tarefa interessante uma vez que nos colocou no lugar de projetistas de uma linguagem e nos confrontamos com problemas e dificuldades que muitas vezes observamos em outras linguagens. Percebemos que, muitas vezes, ao facilitar o desenvolvimento da linguagem para o projetista, o desenvolvedor final que utilizará a linguagem ganha mais responsabilidade e pode ter mais dificuldade com o seu uso. Por sua vez, o contrário também é válido, uma vez que a implementação e estruturação de processos mais complexos por parte dos projetistas pode facilitar a vida do programador final.

Por fim, além da especificação da linguagem, a construção do analisador léxico foi relevante para a fixação do conteúdo próprio da disciplina, apresentando-se como uma tarefa desafiadora, mas interessante.

REFERÊNCIAS

- 1 TAMOSAUSKAS, T. Filosofia chinesa: pensadores chineses de todos os tempos. **Kindle Edition**, 2020.
- 2 LEGGE, J. **The I Ching**. [S.l.]: Courier Corporation, 1963. v. 16.
- 3 MCCRACKEN, D. D.; REILLY, E. D. Backus-naur form (bnf). In: **Encyclopedia of Computer Science**. [S.l.: s.n.], 2003. p. 129–131.
- 4 NIEMANN, T. **A Compact Guide to Lex & Yacc**. [S.l.]: ePaper Press, ????