

TOPGRAPH: uma linguagem para percorrimento em topologias de grafos

Nur Corezzi - 15/0143290

Departamento de Ciência da Computação, Campus Darcy Ribeiro, Asa Norte,
Brasília - DF, 70910-900

Apenas a função de *DFS* não foi implementada por falta de tempo.

1 Introdução

Neste trabalho será apresentada a especificação de uma linguagem a ser desenvolvida e que se destina a área de processamento de grafos. A mesma tem como objetivo fornecer funcionalidades básicas de uma linguagem imperativa como comandos iterativos *for* e *while*, condicionais *if*, tratamento de expressões *booleanas* e aritméticas e também algumas funcionalidades orientadas a grafos. Seu diferencial está na introdução de um tipo básico *graph* que possui uma interface para inserção de arestas direcionais e armazenamento de valores em vértices. Também existirá um comando *for* facilitador que permite iterar fazendo percorrimientos *DFS* (*Depth First Search*) e *BFS* (*Breadth First Search*). A nova primitiva facilita buscas que necessitam ser feitas com frequência em grafos e permite que ações sejam seletivamente executadas de acordo com o percorrimento e os valores dos vértices armazenados. Este tipo básico poupa a necessidade do programador em manter e cuidar do uso de sua própria estrutura, permitindo que algum comportamento seja aplicado de forma prática à topologia de um grafo especificado.

2 Visão geral do tradutor

O tradutor implementado consiste em uma separação de processamento comumente utilizada no desenvolvimento de compiladores para linguagens de programação. As etapas consistem em *análise léxica*, *sintática*, *semântica* e por fim *geração de código intermediário*. Inicialmente será definida uma gramática para a linguagem, em seguida serão determinadas as unidades léxicas para gerar as definições do analisador léxico (gerado em FLEX) que será responsável por fornecer os *tokens* reconhecidos. O analisador sintático criado por *BISON* irá fazer uso dos *tokens* formatados para gerar a árvore de derivação e popular a tabela de símbolos. Por último o analisador semântico irá verificar erros e gerar as anotações devidas na árvore obtida para que em seguida seja feita a geração de código intermediário. Vale ressaltar que apesar da divisão de tarefas a tradução será realizada em apenas uma passada pelo código fonte.

3 Gramática

3.1 Definição

$\langle \text{init} \rangle$	$::= \langle \text{program} \rangle$
$\langle \text{program} \rangle$	$::= \epsilon$ $\langle \text{program} \rangle \langle \text{function} \rangle$
$\langle \text{function} \rangle$	$::= \langle \text{type} \rangle \langle \text{dimension} \rangle \langle \text{id} \rangle ' (' \langle \text{params} \rangle ') ' \langle \text{block} \rangle$
$\langle \text{params} \rangle$	$::= \epsilon$ $\langle \text{params} \rangle ' , ' \langle \text{declaration} \rangle$ $\langle \text{declaration} \rangle$
$\langle \text{function-call} \rangle$	$::= \langle \text{id} \rangle ' (' \langle \text{params-call} \rangle ') '$
$\langle \text{params-call} \rangle$	$::= \epsilon$ $\langle \text{params-call} \rangle ' , ' \langle \text{expr-assign} \rangle$ $\langle \text{expr-assign} \rangle$
$\langle \text{graph-call} \rangle$	$::= \text{'dfs' } \langle \text{graph-params-call} \rangle$ $\text{'bfs' } \langle \text{graph-params-call} \rangle$
$\langle \text{graph-params-call} \rangle$	$::= ' (' \langle \text{id} \rangle ' , ' \langle \text{expr-assign} \rangle ') '$
$\langle \text{statements} \rangle$	$::= \epsilon$ $\langle \text{statements} \rangle \langle \text{statement} \rangle$
$\langle \text{block} \rangle$	$::= \text{'{ ' } \langle \text{statements} \rangle \text{'}' }$
$\langle \text{statement-control} \rangle$	$::= \langle \text{block} \rangle$ $\text{'if' } ' (' \langle \text{expr-assign} \rangle ') ' \langle \text{block} \rangle$ $\text{'if' } ' (' \langle \text{expr-assign} \rangle ') ' \langle \text{block} \rangle \text{'else' } \langle \text{statement-control} \rangle$
$\langle \text{statement} \rangle$	$::= \langle \text{statement-control} \rangle$ $\text{'while' } ' (' \langle \text{expr-assign} \rangle ') ' \langle \text{block} \rangle$ $\text{'for' } ' (' \langle \text{expr-assign} \rangle ' ; ' \langle \text{expr-assign} \rangle ' ; ' \langle \text{expr-assign} \rangle ') ' \langle \text{block} \rangle$ $\text{'for' } ' (' \langle \text{id-or-access} \rangle ' : ' \langle \text{graph-call} \rangle ') ' \langle \text{block} \rangle$ $\text{'>>' } \langle \text{id-or-access} \rangle ' ; '$ $\text{'<<' } \langle \text{expr-assign} \rangle ' ; '$ $\langle \text{declaration} \rangle ' ; '$ $\langle \text{expr-assign} \rangle ' ; '$ $\langle \text{graph-add} \rangle ' ; '$ $\text{'return' } ' ; '$ $\text{'return' } \langle \text{expr-assign} \rangle ' ; '$
$\langle \text{expr-assign} \rangle$	$::= \langle \text{expr-relational} \rangle \text{'=' } \langle \text{expr-assign} \rangle$ $\langle \text{expr-or} \rangle$

$\langle \text{expr-and} \rangle$	$::= \langle \text{expr-or} \rangle \langle \text{and} \rangle \langle \text{expr-and} \rangle$ $\langle \text{expr-or} \rangle$
$\langle \text{expr-or} \rangle$	$::= \langle \text{expr-relational} \rangle \langle \text{or} \rangle \langle \text{expr-or} \rangle$ $\langle \text{expr-relational} \rangle$
$\langle \text{expr-relational} \rangle$	$::= \langle \text{expr-add} \rangle \langle \text{compare-op} \rangle \langle \text{expr-relational} \rangle$ $\langle \text{expr-add} \rangle$
$\langle \text{expr-add} \rangle$	$::= \langle \text{expr-sub} \rangle \langle \text{add} \rangle \langle \text{expr-add} \rangle$ $\langle \text{expr-sub} \rangle$
$\langle \text{expr-sub} \rangle$	$::= \langle \text{expr-mul} \rangle \langle \text{sub} \rangle \langle \text{expr-sub} \rangle$ $\langle \text{expr-mul} \rangle$
$\langle \text{expr-mul} \rangle$	$::= \langle \text{expr-div} \rangle \langle \text{mul} \rangle \langle \text{expr-mul} \rangle$ $\langle \text{expr-div} \rangle$
$\langle \text{expr-div} \rangle$	$::= \langle \text{expr-unary} \rangle \langle \text{div} \rangle \langle \text{expr-div} \rangle$ $\langle \text{expr-unary} \rangle$
$\langle \text{expr-unary} \rangle$	$::= \langle \text{unary} \rangle \langle \text{factor} \rangle$ $\langle \text{factor} \rangle$
$\langle \text{factor} \rangle$	$::= '(\langle \text{expr-assign} \rangle)'$ $\langle \text{value} \rangle$ $\langle \text{function-call} \rangle$
$\langle \text{unary} \rangle$	$::= '!'$ $\langle \text{add} \rangle$ $\langle \text{sub} \rangle$
$\langle \text{and} \rangle$	$::= '&&'$
$\langle \text{or} \rangle$	$::= ' '$
$\langle \text{add} \rangle$	$::= '+'$
$\langle \text{sub} \rangle$	$::= '-'$
$\langle \text{mul} \rangle$	$::= '*'$
$\langle \text{div} \rangle$	$::= '/'$
$\langle \text{compare-op} \rangle$	$::= '<' '<=' '>' '>=' '==' '!='$
$\langle \text{graph-add} \rangle$	$::= \text{'addv' } '(\langle \text{id-or-access} \rangle)'$ $\text{'adda' } '(\langle \text{id-or-access} \rangle, \langle \text{expr-assign} \rangle, \langle \text{expr-assign} \rangle)'$
$\langle \text{declaration} \rangle$	$::= \langle \text{type} \rangle \langle \text{dimension} \rangle \langle \text{id} \rangle$
$\langle \text{dimension} \rangle$	$::= \epsilon$ $\text{'[} \langle \text{number-int} \rangle \text{'}$

$\langle value \rangle$	$::= \langle id-or-access \rangle$ $\langle number \rangle$ $\langle boolean-const \rangle$
$\langle id-or-access \rangle$	$::= \langle id \rangle$ $\langle id \rangle \langle access-lvl \rangle$
$\langle access-lvl \rangle$	$::= \epsilon$ $'[\langle expr-assign \rangle]'$ $\langle access-lvl \rangle$
$\langle type \rangle$	$::= 'int'$ $'boolean'$ $'float'$ $'graph'$ $'void'$
$\langle number \rangle$	$::= \langle number-int \rangle$ $\langle number-float \rangle$
$\langle number-int \rangle$	$::= \langle digit \rangle \langle number-list \rangle$
$\langle number-float \rangle$	$::= \langle digit \rangle \langle number-list \rangle '.' \langle digit \rangle \langle number-list \rangle$
$\langle number-list \rangle$	$::= \langle number-list \rangle \langle digit \rangle$ ϵ
$\langle digit \rangle$	$::= '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$
$\langle boolean-const \rangle$	$::= 'true'$ $'false'$
$\langle id \rangle$	$::= \langle letter \rangle \langle letter-or-digit-list \rangle$
$\langle letter-or-digit-list \rangle$	$::= \langle letter-or-digit-list \rangle \langle letter \rangle$ $\langle letter-or-digit-list \rangle \langle digit \rangle$ $''$
$\langle letter \rangle$	$::= 'A' \mid 'B' \mid 'C' \mid 'D' \mid 'E' \mid 'F' \mid 'G' \mid 'H' \mid 'I' \mid 'J' \mid 'K'$ $'L' \mid 'M' \mid 'N' \mid 'O' \mid 'P' \mid 'Q' \mid 'R' \mid 'S' \mid 'T' \mid 'U'$ $'V' \mid 'W' \mid 'X' \mid 'Y' \mid 'Z' \mid 'a' \mid 'b' \mid 'c' \mid 'd' \mid 'e' \mid 'f'$ $'g' \mid 'h' \mid 'i' \mid 'j' \mid 'k' \mid 'l' \mid 'm' \mid 'n' \mid 'o' \mid 'p' \mid 'q' \mid 'r'$ $'s' \mid 't' \mid 'u' \mid 'v' \mid 'w' \mid 'x' \mid 'y' \mid 'z' \mid '_'$

1. Refatoração de *statements* comandos de *if-else while* e *for* foram simplificados para sempre requisitarem corpo com abertura e fechamento de aspas. A estrutura anterior estava dificultando demasiadamente a geração de código intermediário.
2. Adição de *return* sem parâmetro.
3. Mudança para que expressões relacionais sejam avaliadas antes de *or* e *and*.

4. Remoção do parâmetro de valor a ser filtrado nas chamadas *dfs* e *bfs*. A linguagem permite acesso direto aos valores dos vértices portanto a filtragem automática não faz muito sentido, porém o percorrimento em topologia ainda se faz útil.

3.2 Aspectos informais da gramática

A gramática em questão define uma lista de funções como um programa válido, tendo ponto de execução inicial em uma função definida com identificador *main*. Cada função possui um conjunto de argumentos que são declarações que envolvem um tipo e um identificador e possui também um tipo de retorno. Serão também delimitadas por um conjunto de abre e fecha chaves `{ }` que possui uma lista de declarações.

Uma declaração constitui um comando da linguagem que podem ser qualquer um dos especificados em *statements*. Existe também um comando diferenciado de *for* para o percorrimento dos grafos que contempla também chamadas de funções *BFS* e *DFS* na regra *graph-op*. Esta chamada recebe um identificador do grafo e um parâmetro indicando um vértice de início. A função irá retornar os vértices na ordem do percorrimento especificado (exemplo em código mais abaixo).

Um grafo irá iniciar com zero vértices e possui duas funções auxiliares para acrescentar vértices(*addv*) e arestas(*adda*) de forma dinâmica. Cada vértice armazena um valor inteiro que pode ser acessado por colchetes e será sempre inicializado em zero. Esta decisão evita tipos genéricos na definição apenas para simplificar a implementação. Com relação às operações (descritas mais abaixo em código) será possível atribuir valores a um vértice, identificado por um inteiro, e designar arestas entre vértices existentes.

Os comandos de *for*, *while* e *if* contemplam expressões que irão avaliar os fluxos de execução e também uma regra *statement* que permite encadear uma lista de comandos à sua execução.

Expressões podem ser consideradas como operações aritméticas, lógicas e atribuições, que seguem uma precedência nessa mesma ordem. Existe também a possibilidade de definir prioridades distintas a partir do uso de parênteses definidos na regra *factor*.

Os tipos básicos da linguagem contemplam *boolean*, *int*, *float*, *void* e *graph*. O tipo *graph* é um tipo especial que também é unicamente identificado pelo tamanho especificado entre colchetes *[number-int]*. Todos os tipos podem também ser declarados como *array* unidimensional (indexados a partir de zero) com exceção de *void* que apenas pode estar presente em declaração de funções.

Por último os numerais da linguagem são definidos pelos terminais numéricos como listas que definem inteiros e reais e os identificadores são definidos pelos caracteres alfabéticos minúsculos e maiúsculos além do *underscore*.

3.3 Descrição semântica

A seguir segue um pequeno trecho de código explicitando o uso da nova primitiva em um programa básico da linguagem:

```
void main() {
    graph g;
    // Acrescenta um vertice
    addv(g);
    addv(g);
    //Define uma aresta direcional
    adda(g,0,1)
    // Designa valor para algum vertice
    g[0] = g[1] = 1;
    int vertice;
    for (vertice : bfs(g, 0)) {
        // Percorre o grafo em bfs iniciando no vertice zero
        << vertice;
    }
}
```

O ponto inicial de execução será sempre a função que possui o identificador de *main* e sua falta implica em não execução e será gerado um *warning*. As funções devem ser declaradas antes de serem utilizadas assim como as variáveis. Parâmetros de tipos básicos *boolean*, *int* e *float* serão passados por cópia porém *arrays* e *graph* serão passados por referência. Apenas será possível passar um *array* se tanto o destino como a fonte tiverem o mesmo tipo, o que é estabelecido pelo tamanho, dimensões e tipo armazenado. Não será possível retornar *array* nem *graph* de funções, apenas os tipos *boolean*, *int*, *float* e *void*.

Arrays terão seu tamanho estabelecido em tempo de compilação e apontarão para a mesma região de memória ao longo da execução do programa não podendo ser realocados ou apontar para outro *array*. O tipo *graph* também irá se comportar de forma análoga sem poder participar de atribuições pelo usuário porém será uma estrutura dinâmica em que suas dimensões podem crescer.

Com relação a conversões implícitas a única necessária e existente será entre *float* e *int*. O tipo *boolean* é implementado como um *int* onde qualquer valor diferente de zero será *true* e caso contrário *false*. Caso exista um *float* na operação a *cast* dos operandos será para *float*, e caso destino seja *int* ou *boolean* será efetuado um *cast* de *float* para *int*. Portanto na árvore abstrata apenas existem anotações de conversões *floatToInt* e *intToFloat*. Vale ressaltar também que expressões de atribuição possuem um valor de retorno que será o do próprio identificador sendo atribuído

O escopo das variáveis é delimitado pelas chaves. A busca de uma variável na tabela se resume a achar a ocorrência do símbolo desejado onde o escopo atual está contido no escopo do símbolo encontrado. Caso haja diversos símbolos na hierarquia que sejam candidatos então o símbolo do escopo mais próximo acessível será escolhido. Antes de ser inserida uma variável a mesma busca é efetuada para verificar redeclarações. A busca também se aplica para funções

com uma restrição adicional de que as mesmas são identificadas pelo seu nome de declaração portanto não podem existir funções distintas com o mesmo nome, mesmo que a lista de parâmetros seja diferente.

Os comandos condicionais e iterativos exigem que a declaração de variáveis seja feita fora de sua assinatura e possuem um escopo próprio (com exceção de *if* simples sem chaves). O iterador de grafos *for* exige que o parâmetro recebendo o vértice de retorno seja um *int* uma vez que os vértices apenas são identificados como inteiros.

4 Implementação

Para a realização da análise léxica foi utilizado o gerador de analisador léxico *FLEX*. Foi realizada a definição regular para reconhecimento dos lexemas para que em seguida fosse transportada para as definições em *FLEX*. Os tokens reconhecidos são repassados como uma enumeração, a informação adicional sobre o lexema é encontrada nas próprias variáveis de *buffer* que o gerador produz podendo ser acessadas pelo analisador sintático com facilidade.

O gerador de analisador sintático escolhido foi o BISON. Este gerador cria analisadores *bottom-up* LALR(1) e facilita o procedimento de criação dos conjuntos e estados LR, gerando também as tabelas de ações e os procedimentos necessários para a derivação da árvore sintática. Seguindo a especificação da documentação, a gramática formalmente apresentada acima foi traduzida para o formato de produções compreendidas pelo BISON. As definições de tokens e regras da gramática podem ser encontradas no arquivo fonte *gramatica.y*.

Para cada regra da gramática foram associadas ações semânticas para a criação da árvore abstrata e análise semântica. Cada nó da árvore é representado por uma estrutura genérica *Node* que irá conter um identificador do símbolo da gramática o qual representa para auxiliar em ações posteriores. Foram abstraídas rotinas adicionais para criação destes nós, assim como para a criação e manipulação da tabela de símbolos que serão melhor especificadas abaixo.

4.1 Definição regular

Na Figura 1 é possível identificar as variáveis em negrito/italico e os símbolos terminais em fonte normal. Não foram agrupados símbolos por funcionalidade como comparadores, operadores, tipos e sim um token por unidade terminal da linguagem. Caso seja necessário agrupamentos para facilitar a análise sintática, posteriormente podem ser definidas funções que determinam pertinência de determinado token em algum grupo. Cada token será identificado por um valor de *enum* gerado pelo analisador sintático. Os tokens são especificados em *gramatica.y* pelas regras prefixadas em *%token*.

4.2 Estrutura da árvore sintática abstrata (AST)

A árvore gerada é construída por meio de uma abstração de um nó. Cada nó pode conter diversos filhos e possuem um identificador para informar qual símbolo da

ws → [\t\n]	read → >>
letter → [A-Za-z]	write → <<
digit → [0-9]	and → &&
	or →
if → if	le → <=
else → else	ge → >=
for → for	less → <
while → while	greater → >
	eq → ==
boolean → boolean	ne → !=
int → int	not → !
float → float	mul → *
void → void	div/ → /
graph → graph	sum → +
to → to	sub → -
	assign → =
true → true	end → ;
false → false	open_brace → {
	close_brace → }
dfs → dfs	it → :
bfs → bfs	separator → ,
return → return	open_p → (
	close_p →)
id → (letter) (letter digit _)*	open_bracket → [
number → (digit)+ (. digit)+?	close_bracket →]

Figura 1. Definições regulares, **number** foi separado em **c_int** e **c_float** para facilitar identificação.

gramática o gerou. Os nós são criados individualmente a cada redução e a estrutura da árvore é gerada a partir do método de adição de filhos. Cada regra composta por mais de um símbolo, que já possua valor semântico presente (ou seja nó instanciado), será acrescida como filho do símbolo ao qual a regra será reduzida. Derivações que geram listas terão seu formato em árvore simplificado para listas ligadas, o que facilita alguns processamentos. Também houve a remoção de nós não essenciais da sintaxe após as reduções, como vírgulas, parênteses e etc. Ao final, na regra geradora o valor será atribuído a um ponteiro *AST* declarado no analisador.

Cada nó possui uma informação de qual símbolo da gramática o mesmo representa como por exemplo uma função, parâmetro, tipo, identificador e etc. Os nós também possuem ponteiros para seus filhos e além disso uma informação adicional de complemento utilizada por variáveis, funções e constantes numéricas que precisam do valor ao qual se referem. Nós também possuem uma *TypeExpression* que define seu tipo na gramática. O tipo será utilizado para verificações de operações possíveis e conversões em *casts* implícitos. As conversões também serão anotadas na árvore especificamente no nó com um campo informando se a conversão deve ser *floatToInt* ou *intToFloat*, que por enquanto são as duas únicas conversões possíveis. Vale ressaltar que por simplicidade os tipos *boolean* são traduzidos e tratados como tipos *int* pelo tradutor. **Por último o código gerado em TAC para determinado nó será anexado como informação e uma lista de *back-patch* também estará presente para ser resolvida a medida que a lista de *statements* for sendo derivada.**

4.3 Tabela de símbolos

A tabela de símbolos inicialmente é uma estrutura de lista de listas ligada onde cada colisão gera um novo encadeamento em alguma entrada. Os elementos instanciados são declarações de variáveis e de funções. As informações adicionadas são elementos de uma enumeração indicando se é uma função ou variável, um identificador com o nome da variável ou função, o tipo do símbolo na gramática (*float*, *int*, *graph<int>*, *arrays* e etc), que irá facilitar verificações e uma referência para o nó da árvore abstrata caso seja necessário ter acesso a informações adicionais, como por exemplo os parâmetros de funções para verificações semânticas. **O símbolo em TAC que guarda a região de memória para o objeto alocado também estará presente na tabela para futuras referências.**

Também foi adicionada uma informação adicional de escopo. Cada entrada da tabela possui um escopo próprio que é construído juntamente com a árvore abstrata. A cada abertura de chaves um novo escopo é adicionado na pilha de escopo global e são representados por um inteiro obtido de um contador. O escopo corrente em qualquer momento será a concatenação de todos os escopos nesta pilha e será também a informação adicionada juntamente com os símbolos na tabela. Esta informação será utilizada na busca de símbolos na tabela para que um símbolo seja referenciado de forma que faça sentido seguindo a estrutura da árvore abstrata construída.

4.4 Gerador de código intermediário

A tradução para *Three Address Code*(TAC) será realizada em uma passagem, portanto será necessário a realização de *backpatching* para a resolução de pulos para rótulos ainda não definidos. Cada nó possui uma lista *next* que irá conter os campos que devem ser atualizados com algum rótulo futuro. Expressões serão avaliadas por completo evitando a realização de *short-circuit* para simplificação das rotinas de geração de expressões. Cada expressão gerada na linguagem será atribuída a uma variável temporária em TAC e poderá ser acessada a partir do código que irá ser armazenado em cada nó da árvore abstrata.

Cada início de função reinicia o contador de variáveis temporárias para que possam ser reutilizadas no novo contexto. Devido a falta de alocações estáticas globais, todas as declarações na linguagem irão gerar um endereço armazenado em variável temporária e a sua referência será armazenada na tabela de símbolos para posteriores acessos e liberação de alocações. Já para parâmetros de funções seus valores serão armazenados em pilha e o símbolo na tabela não será um temporário e sim uma referência para a variável em pilha.

Para *arrays* serão armazenados no início do endereço do código em TAC o tamanho fazendo com que exista um *offset* para os reais valores armazenados. Endereços em TAC sempre serão acompanhados de um índice para seu acesso, uma vez que operações aritméticas com endereços não devem ser realizadas.

Variáveis do tipo *graph* irão alocar inicialmente memória para um valor caracterizando o tamanho do grafo, outro ponteiro para o *array* de valores dos vértices e um terceiro ponteiro para um *array* de arestas do grafo. Tanto o *array* de vértices quanto o de arestas serão alocados com tamanho inicial em zero. A lista de arestas é um *array* de ponteiros para listas uma vez que cada lista de arestas pode conter um tamanho variado e portanto não serão armazenados de forma contígua. Cada adição de aresta ou vértice gera novas alocações e atualizam os valores nos campos armazenados na tabela de símbolos.

Para facilitar os procedimentos de geração de código foi gerada uma estrutura que define instruções de código em TAC assim como funções para manipular a geração. Cada *Instruction* possui um rótulo, o comando em TAC o qual representa e os parâmetros representados por três campos *Field*. Instruções que não utilizam todos os campos ou rótulos apenas devem passar valores nulos nos campos não desejados. Também existe um campo *next* em instruções que corresponde a próxima instrução encadeada. Funções auxiliares foram geradas para facilitar a criação de campos, rótulos, instruções, junção de instruções e geração de estruturas mais complexas na linguagem como por exemplo declarações ou rotinas de *output* e *input* as quais serão descritas mais abaixo.

4.5 Descrição de funções adicionais

Funções adicionais foram criadas para: auxiliar na criação da tabela de símbolos, gerar os nós da árvore abstrata, gerar código intermediário, fazer liberação de alocações e auxiliar na impressão das estruturas para a saída. A seguir se encontram algumas das funções adicionais utilizadas pelo programa tradutor:

Árvore sintática abstrata

1. *create_node* e *push_child*: Responsáveis por criar um nó *standalone* e adicionar filhos a algum nó.
2. *free_node*, *free_tree*: Responsáveis por liberação de memória alocada para a árvore abstrata.
3. *print_node_prefix*, *print_node_suffix*, *print_node* e *print_tree*: Funções auxiliares para imprimir no console a árvore abstrata formatando os prefixos e separando nós de forma individual para depuração.

Tabela de símbolos

1. *stable_create_symbol*, *stable_find* e *stable_add*: Funções auxiliares para criar um símbolo *standalone*, identificar a existência de um símbolo na tabela e adicionar um símbolo qualquer a tabela.
2. *stable_find_with_scope*: Realiza uma busca na tabela levando em consideração um escopo atual. Será utilizada em conjunto com um ranqueamento de acesso a escopos para determinar símbolo mais próximo do escopo atual.
3. *free_symbol_node* e *free_stable*: Responsáveis por liberação de memória alocada para a tabela de símbolos.
4. *stable_symbol_print* e *stable_print*: São funções auxiliares para imprimir no console símbolos individuais e tabelas de símbolos.

Escopo

1. *scope_access_rank*: Estabelece uma classificação de acessibilidade de um escopo por outro. Irá informar se um escopo B pode acessar um escopo A e se é acessível qual o nível de diferença entre os mesmos. Função irá auxiliar nas buscas da tabela de símbolos.
2. *scope_push* e *scope_pop*: Irá atualizar uma referência global de escopo em formato de pilha, *scope_push* em particular recebe uma referência da função atual na tabela de símbolos para facilitar alguns processamentos.
3. *scope_sz*, *scope_eq*, *scope_create*, *cpy_scope*, *print_scope* entre outras: São funções auto-explicativas que fornecem auxílio no tratamento de escopos.

Expressões de Tipos

1. *type_get_cast*: Retorna o tipo de conversão que será necessário entre um tipo de origem e um tipo de destino.
2. *type_max*: Em uma hierarquia de tipos retorna o tipo comum mais próximo entre os dois tipos passados como parâmetros.
3. *type_from_access_lvl*: Devido a acesso de variáveis por colchetes é necessário inferir o tipo sendo tratado no acesso.
4. *type_can_assign*, *type_can_return*, *type_is_arithmetic* e *type_is_boolean*: Informa como alguns tipos podem ser tratados, se podem estar em uma atribuição, retorno de função, em uma expressão aritmética ou expressão booleana.

5. *type_build*, *type_cpy*, *type_eq*, *free_type*, *type_print* entre outras: São funções também auto-explicativas e fornecem abstrações para as funções acima e funcionalidades gerais ao tradutor.

Gerador de código intermediário

1. *cgen_label*, *cgen_field* e *cgen_instruction*: Funções auxiliares para gerar as instruções em TAC.
2. *cgen_patch*, *cgen_cpy_patch*, *cgen_merge_patch* e *cgen_back_patch*: Auxiliam no tratamento e resolução de *labels*.
3. *cgen_declaration* e *cgen_declaration_param*: geram código de declaração de variáveis em parâmetros e dentro de funções, também populam a tabela de símbolos com referência para os campos que guardam os endereços alocados.
4. *cgen_function_call* e *cgen_declaration_param*: Geram código para declaração de variáveis em parametro de função e de chamadas de função.
5. *cgen_var_access*: Gera código para acesso a valor de variáveis, resolve casos em que o que será retornado será um *lvalue* ou um *rvalue*.
6. *cgen_expression_relational*, *cgen_expression_boolean* e *cgen_expression_arithmetic*: Geram código para as expressões da linguagem.
7. *cgen_if* e *cgen_if_else*: Geram código para as estruturas de controle.
8. *cgen_write* e *cgen_write_array*: Gera instruções de escrita em console para variáveis com *lvalue* e *rvalue*.
9. *print_instruction* e *print_code*: Geram código TAC em saída padrão.

5 Tratamento de erros

5.1 Identificação e tratamento de erros léxicos

A nível de análise léxica apenas será necessário verificar erros referentes a caracteres não existentes e padrões que contemplem símbolos do alfabeto da linguagem mas que não se encaixem nas definições regulares informadas. A forma de resposta será sempre informar caractere a caractere do padrão não reconhecido informando linha e coluna. Em seguida todos os caracteres incorretos serão descartados da análise e o analisador irá continuar com o restante da entrada. Não foi necessário a introdução de procedimentos para tratamento de erro uma vez que a regra de captura de padrões não reconhecidos (regra ". {}" em FLEX) irá automaticamente descartar os caracteres não desejados bastando apenas acrescentar uma rotina para informar a ocorrência do erro.

5.2 Identificação e tratamento de erros sintáticos

Por simplicidade nenhum erro identificado será corrigido. O *token* responsável pelo erro será simplesmente descartado e o analisador por padrão irá desempilhar os estados até encontrar alguma derivação que possua um *token error*. Este *token* se apresenta em pontos chaves do código sendo eles a regra de produção de um

função e a regra de produção de *statements* para que a partir destes pontos o analisador possa retomar apenas para verificar possíveis próximos erros na derivação. Os erros serão descritos identificando a linha e coluna do *token* o qual gerou o erro, e será informado qual *token* deveria ser esperado, sempre que possível inferir.

5.3 Identificação e tratamento de erros semânticos

Um conjunto não exaustivo de erros semânticos foi determinado levando em consideração o tempo disponível para implementação. Os erros são informados com a mensagem respectiva a linha e coluna do lexema que pode ter gerado o erro. A seguir seguem os erros:

1. Variável ou função já declarada, caso já apresentem uma ocorrência de mesmo escopo na tabela de símbolos.
2. Variável ou função não encontrada no escopo, caso não exista ocorrência na tabela de símbolos.
3. Operandos incompatíveis em expressões aritméticas e booleanas como por exemplo operações entre *arrays* e grafos que não existem na linguagem.
4. Tipos a esquerda em expressões de atribuição apenas podem ser *lvalue*.
5. Chamada indevida caso os tipos instanciados para a chamada de função não possuam uma equivalência de tipos.
6. Incompatibilidade do tipo em *return* e tipo da função.
7. Incompatibilidade de tipos em operações de atribuição.
8. *void* não pode ser tipo em variável nem array.
9. Funções não podem retornar *array* nem *graph*.
10. Tamanho da dimensão de *arrays* na declaração não pode ser zero.
11. Acesso a elemento de *array* deve ocorrer apenas com índice inteiro.

6 Arquivos de teste

6.1 Arquivos de teste léxico

Os arquivos *lexico/correto1*, *lexico/correto2*, *lexico/errado1* e *lexico/errado2* possuem exemplos de programas sem e com erros léxicos. Em *errado1* são apresentados caracteres não reconhecidos pela linguagem, espera-se que sejam informados um a um os caracteres incorretos, e em *errado2* erros com símbolos do alfabeto da linguagem, porém com padrões não reconhecidos onde será apresentado um a um os caracteres que não fazem parte de um padrão reconhecível. Ambos arquivos de erro são descritos na Figura 4, lembrando que os comentários apenas descrevem o erro e não necessariamente o que será informado ao usuário.

```

void main() {
    @""\çã~^ // Cada simbolo será informado como erro
}

void main() {
    .9 1.      // Numeros incorretos
    a | b      // Apenas um pipe não reconhecido
    a &&& b     // & a mais
    _nome      // Underscore no inicio de identificador
}

```

Figura 2. Programas comentados com os respectivos erros léxicos

6.2 Arquivos de teste sintáticos

Os arquivos *sintatico/correto1*, *sintatico/correto2*, *sintatico/errado1* e *sintatico/errado2* possuem exemplos de programas sem e com erros sintáticos. Em *errado1* é apresentada uma função sem tipo de retorno declarado, um for sem declaração de uma variável na primeira parte e uma declaração de variável sem fechamento de ponto e vírgula. Em *errado2* existem erros de comando de leitura sem variável de destino, for iterador sem uma função correta de iteração, um if sem condicional e um acesso incorreto de um vértice por meio de ponto flutuante. Em cada caso espera-se que o token imediato do erro seja informado e em algumas situações, onde seria possível inferir o token esperado, o mesmo também será informado. Em alguns casos é possível inferir qual seria o token esperado mas em muitas situações existe ambiguidade no que seria a intenção do usuário, portanto algumas mensagens não são precisas e difíceis de se determinar.

```

1  main() {           // Função sem tipo será identificado como ID não esperado
2      int a, b, c;   // Todos as construções já reduzidas que compõem a função serão descartadas
3  }
4
5  void main () {
6      int a = 2      // Seria esperado um ';' porém foi encontrado o token FOR
7      for (; i; i) {} // For sem declaração de índice para iterar, não é possível determinar um erro exato
8                      // pois neste caso diferentes construções corretas poderiam ser sugeridas
9      while(a) a;
10 }

1  void main() {
2      >>;            // ';' não esperado seria esperado algum valor de uma expressão como ID
3      for (int i : vdsf(a, 1, 2)) { // FOR de iteração deve conter alguma função DFS ou BFS
4      }
5      if () {}       // IF sem expressão para avaliar, final de parentesis não esperado
6                      // Acesso a vértice deve ser por inteiros, porém aqui mais de uma
7      g[1.2];        // construção será possível portanto mensagem talvez não seja significativa
8
9  }

```

Figura 3. Programas comentados com os respectivos erros sintáticos

6.3 Arquivos de teste semânticos

Os arquivos *semantico/correto1*, *semantico/correto2*, *semantico/errado1* e *semantico/errado2* possuem exemplos de programas sem e com erros semânticos. Os erros esperados foram descritos nas figuras abaixo com comentários ao lado da linha em que se espera que o erro ocorra.

```
1 void main() { }
2
3 void main () {      // Função redeclarada
4     int a;
5     a = 2;
6     int a;          // Variável redeclarada
7     {
8         int a;
9         outra();     // Função não encontrada no escopo
10        b = 2;       // Variável não encontrada no escopo/tipos incompatíveis
11        graph c;
12        c = a;       // Tipos incompatíveis em atribuição
13        l = 1;       // Variável a esquerda em atribuição de ser lvalue
14    }
15    int b;
16 }
```

```
1 graph outra(int a) { // graph não pode ser retorno
2     return a;        // Tipo incompatível com retorno
3 }
4
5 void main() {
6     graph g;
7     int a;
8     a = 1 + g;       // Tipos incompatíveis em expressão e atribuição
9     outra(g);        // Chamada com parâmetros incompatíveis
10    outra();          // Chamada com parâmetros incompatíveis
11 }
```

Figura 4. Programas comentados com os respectivos erros sintáticos

7 Dificuldades da implementação

Com relação a parte léxica as definições da linguagem se mostraram simples e não exigiram grande aprofundamento no ferramental fornecido pelo *FLEX*. Na parte sintática já foi necessário maior estudo com relação as funcionalidades fornecidas pelo *BISON*. Ambiguidades presentes na gramática também apresentaram desafios na sua redefinição. Também existiram dificuldades com relação ao que seria necessário para a próxima fase uma vez que não havia muito conhecimento sobre o que seria feito na análise semântica nem na geração de código intermediário. Essa dificuldade já se apresentou na parte léxica uma vez que mudanças foram efetuadas nesta fase devido ao desconhecimento de certas necessidades.

Na parte semântica a maior dificuldade foi determinar um subconjunto de erros e conversões que fossem possíveis ser tratados no tempo disponível para implementação. A verificação da existência de um retorno em funções foi descartada pois exige a verificação de todos os caminhos possíveis no código de uma função o que é dificultado pelos comandos condicionais. Também foram descartados *cast* entre *arrays* uma vez que podem se apresentar de forma complexa na geração de código intermediário. Também foram realizadas algumas simplificações na linguagem para poder generalizar melhor os tratamentos sendo efetuados nesta fase.

Para gerar o código intermediário a maior dificuldade encontrada foi tratar endereços em TAC. Uma abstração inicial implementada seguia a manipulação direta de endereços realizando *offset* por somas para manipular diferentes endereços porém a documentação do TAC não recomenda que seja feito desta forma e sim pela utilização das funções disponíveis que exigem um índice para calcular o endereço. Portanto uma refatoração significativa teve de ser realizada para anexar um campo de índice a qualquer endereço sendo manipulado.

Referências

- [ALSU07] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 2 edition, jan 2007.
- [Dega] Jutta Degener. ANSI C grammar, lex specification. <http://www.quut.com/c/ANSI-C-grammar-l-2011.html>. Accessed: 2020-03-17.
- [Degb] Jutta Degener. ANSI C yacc grammar. <http://www.quut.com/c/ANSI-C-grammar-y-2011.html>. Accessed: 2020-03-17.