



Volume 1

5
SEMINÁRIO INTEGRADO
SOFTWARE E HARDWARE
NACIONAIS

UNIVERSIDADE FEDERAL DO RIO DE JANEIRO
Núcleo de Computação Eletrônica
CAPRE/FINEP

1978

004.06

S471

1978 V.1

GALE - Gerador de Filtros Lêxicos

Conferencista e Autor:
Arndt von Staa
PUC/RJ, Informática

Resumo:

No processo de compilação a primeira atividade a ser exercida é a de decompor o texto fonte em unidades lêxicas. Durante esta decomposição, serão encontradas unidades lêxicas que não possuem utilidade do ponto de vista do processamento posterior, tais como comentários e sequências de "brancos". Estas unidades "inúteis" devem ser eliminadas pelo filtro lêxico.

Por sua vez, as unidades lêxicas úteis poderão ser símbolos representados por meio de cadeias de caracteres de tamanho variável, por exemplo identificadores de variáveis. O processamento de cadeias de tamanho variável é custoso. É desejável, então, utilizar-se uma identificação de tamanho padronizado e curto, em correspondência um para um com este símbolo. Esta compressão é efetuada por uma tabela de símbolos.

O filtro lêxico é, portanto, uma rotina capaz de decompor o texto fonte em unidades lêxicas, eliminar aquelas que forem inúteis, comprimir as que forem úteis através de uma tabela de símbolos, retornando identificações padronizadas, de tamanho curto e constante.

O gerador de filtros lêxicos é um processador de linguagem, que traduz uma especificação de um filtro lêxico dada em linguagem própria para um programa capaz de perfazer a filtragem lêxica tal como especificada.

O presente trabalho apresenta os aspectos considerados mais relevantes da especificação lógica do sistema gerador de filtros lêxicos. A especificação completa encontra-se em [Staa78].

1 Introdução.

No processo de compilação a primeira atividade a ser exercida é a de decompor o texto fonte em unidades léxicas. Esta atividade pode tanto ocorrer como fase isolada no processo de compilação, quanto como subatividade de um processamento mais complexo.

Quando da especificação de um sistema gerador, deverá ser imposto um mínimo de restrições, uma vez que é virtualmente impossível prever-se todas as aplicações possíveis. Consequentemente, necessita-se de interação com o usuário para fins de leitura, de impressão do texto fonte e, possivelmente, de outras tarefas.

Devido à necessidade de minimizar restrições, o sistema torna-se mais complexo, uma vez que passa a ter mais opções. Esta complexidade maior não deve incidir sobre filtros específicos criados para situações particulares. Por esta razão, optou-se por um sistema gerador produzindo não somente tabelas, mas também programas, pois, desta forma, é possível eliminar-se código irrelevante de filtros léxicos específicos.

Neste manual especificaremos primeiro como deverão funcionar os filtros léxicos e, depois, como definir filtros léxicos específicos.

2 Especificação do Filtro Léxico.

2.1 Definição.

Um filtro léxico é uma rotina que, a cada ativação, reconhece uma ou mais unidades léxicas, das quais a última foi definida como sendo útil para fases posteriores e todas as unidades léxicas anteriores à última foram definidas como sendo inúteis para fases posteriores e, portanto, deverão ser eliminadas durante a filtragem léxica.

A cada ativação o filtro léxico retorna uma identificação em correspondência um para um com a cadeia de caracteres que constitui a unidade léxica útil lida. Esta identificação, que chamaremos de lexema, possui tamanho padrão e é comprimida através de uma tabela de símbolos[DeRemer74].

São exemplos de unidades léxicas:

operadores, por exemplo: + - * / , .
variáveis
constantes
comentários- estes são eliminados durante a
filtragem léxica e não são lexemas

Os filtros l\u00e9xicos s\u00e3o gerados pelo sistema gerador de filtros l\u00e9xicos a partir de uma linguagem de defini\u00e7\u00e3o pr\u00f3pria. Cada filtro l\u00e9xico consta de tres componentes principais:

- i- rotinas de controle e inicializa\u00e7\u00e3o;
- ii- transdutor finito reconhecedor de unidades l\u00e9xicas que ser\u00e1 chamado de analisador l\u00e9xico;
- iii- tabela de s\u00edmbolos.

Os filtros l\u00e9xicos permitem:

- i- decompor o texto fonte em unidades l\u00e9xicas conforme defini\u00e7\u00e3o dada pelo usu\u00e1rio;
- ii- descartar unidades l\u00e9xicas definidas como in\u00fateis para fases posteriores \u00e0 filtragem l\u00e9xica;
- iii- retornar os lexemas correspondentes a unidades l\u00e9xicas consideradas \u00fateis para fases posteriores;
- iv- ativar rotinas externas sempre que for assim determinado pelo analisador l\u00e9xico definido pelo usu\u00e1rio;
- v- descartar e/ou substituir e/ou incluir caracteres, ou seq\u00fancias de caracteres, no interior dos s\u00edmbolos correspondentes a lexemas;
- vi- reconhecer palavras chave reservadas e/ou pr\u00e9-definidas;
- vii- perfazer o reconhecimento condicional de se\u00e7\u00f5es componentes de uma unidade l\u00e9xica;
- viii- permitir o acesso a vari\u00e1veis internas atrav\u00e9s de par\u00e2metros passados para rotinas do usu\u00e1rio.

Os filtros l\u00e9xicos n\u00e3o far\u00e3o:

- i- impress\u00e3o dos dados lidos. Isto somente poder\u00e1 ser feito atrav\u00e9s de rotinas fornecidas pelo usu\u00e1rio;
- ii- leitura do texto fonte. Isto somente poder\u00e1 ser feito atrav\u00e9s de rotina fornecida pelo usu\u00e1rio. Esta ser\u00e1 ativada por comando efetuado pelo analisador l\u00e9xico. Este comando estar\u00e1 usualmente associado ao reconhecimento da condi\u00e7\u00e3o de fim de texto na \u00e1rea contendo o texto a ser analisado l\u00e9xicamente.

2.2 Entrada.

Os filtros l xicos n o s o capazes de perfazer a leitura do texto fonte. Tal opera o   efetuada atrav s de uma rotina fornecida pelo usu rio. Esta rotina ter  seu nome definido na se o <comunica o> da defini o do filtro l xico.

Sempre que a rotina de leitura for chamada, dever  ser devolvido por ela um ponteiro (endereço) para uma  rea que cont m o texto fonte (ou parte dele) a ser analisado l xicamente.

  seguinte o formato da  rea de texto a ser analisada l xicamente:

```
-----  
| TAM | POS | c1 | c2 |   ...   | cTAM |  
-----
```

TAM- indica o tamanho da cadeia de caracteres contida na  rea de texto a ser analisado l xicamente ("buffer" de entrada). Este valor n o   alterado pelo filtro l xico.

POS- cont m o  ndice do pr ximo caracter a ser lido. Este valor deve ser inicializado sempre pela rotina de leitura fornecida pelo usu rio. A inicializa o dever  ser tal que POS indique para o primeiro caracter da por o da cadeia a ser analisada l xicamente. O valor inicial de POS dever  estar contido no intervalo $1 \leq POS \leq TAM$. O valor de POS   continuamente alterado pelo filtro l xico.

ci- onde $1 \leq i \leq TAM$ s o os caracteres a serem lidos. Estes caracteres poder o ser alterados pelo filtro l xico.

Condi es especiais:

- i- se $TAM \leq 0$,   reconhecido o caracter especial \$FARQ - fim de arquivo.
- ii- se $POS > TAM$,   reconhecido o caracter especial \$FBUF - fim de "buffer". Ao ser encontrado este caracter, a rotina de leitura fornecida pelo usu rio ser  acionada, desde que esta opera o esteja associada ao reconhecimento do caracter especial \$FBUF.
- iii- se a rotina de leitura fornecida pelo usu rio retornar o ponteiro nulo, ser  reconhecido o caracter \$FIM. Uma vez tendo sido reconhecido \$FIM, o filtro l xico passar  a reconhecer somente este caracter especial. Isto ocorrer  independente de chamadas de rotinas fornecidas pelo usu rio e associadas ao reconhecimento de \$FIM. O reconhecimento do caracter \$FIM causa, tamb m, a transfer ncia para arquivo do conte do da tabela de s mbolos neste momento residente na mem ria principal.

iv- diversas condições poderão levar o filtro léxico a reconhecer o caracter especial \$ERRO. São seguintes as condições:

- a- o transdutor léxico é incapaz de reconhecer pelo menos uma unidade léxica, útil ou não;
- b- o transdutor léxico atingiu um estado definido como inconsistente pelo usuário. Associado a este estado existe uma mensagem de erro criada pelo usuário;
- c- o transdutor léxico atingiu um estado definido como inconsistente pelo próprio sistema gerador de filtros léxicos. Associado a este estado existe uma mensagem de erro criada pelo sistema. Exemplos de inconsistências deste gênero são: esgotamento da capacidade de armazenagem de símbolos, tamanho de símbolo excessivamente grande etc.

Em todos estes casos de reconhecimento de condição de erro, o reconhecimento em curso será suspenso e será gerado um lexema indicativo de erro. Em alguns casos, será possível recomençar a operação do filtro léxico imediatamente após ao ponto em que foi encontrada a condição de erro.

v- se e somente se a expressão condicional contida em um reconhecimento condicional for verdadeira, é reconhecido o caracter especial \$COND. Caso esta condição seja falsa, nada será reconhecido por esta expressão. A qualquer instante pode ser reconhecido o "caracter" vazio. Isto corresponde ao reconhecimento de \$COND onde a condição de reconhecimento é sempre verdadeira.

Cabe ao usuário evitar tentativas de leitura além do final de arquivos de texto fonte, através da especificação adequada dos filtros léxicos e/ou através do controle de ativações da rotina de leitura fornecida.

A sistemática de leitura será tal que o usuário do gerador do filtro léxico tenha a impressão de estar lendo uma cadeia "infinita" de caracteres do conjunto:

ALFABETO U {\$FARQ, \$FBUF, \$FIM, \$ERRO, \$COND}

A rotina de leitura fornecida pelo usuário deverá ser capaz de selecionar arquivos e de salvar registros ainda não completamente analisados caso seja permitida a leitura de texto fonte a partir de diversos arquivos e/ou áreas contendo texto fonte.

Os caracteres especiais \$FARQ, \$FBUF, \$FIM, \$ERRO e \$COND não possuem representação binária própria. Sendo assim, também não podem fazer parte da cadeia de caracteres que constitui um símbolo na tabela de símbolos. Estes caracteres cor-

respondem a condições lógicas reconhecidas pelo filtro léxico. Estas condições muitas vezes não são repetíveis, uma vez que dependem diretamente das condições de contorno no momento em que foram reconhecidas.

2.3 Resultados Produzidos.

Para cada ativação do filtro léxico será retornado um descritor de lexema, possuindo o seguinte formato:

```
-----  
| POS | TIPO | ID |  
-----
```

POS- é o índice do primeiro caracter na primeira área de entrada, utilizado para o reconhecimento da unidade léxica. Este índice independe da possibilidade de descarte de caracteres que constituem a unidade léxica lida, não sendo, portanto, transcrito para o símbolo correspondente ao lexema. Por exemplo, quando for reconhecida uma constante "cadeia de caracteres" obedecendo as convenções do PL/1, ambas as aspas inicial e final devem ser descartadas do símbolo antes de pô-lo na tabela de símbolos, o índice POS indica porém, o local onde estava a aspa inicial da cadeia. O índice POS é particularmente útil para a geração de mensagens de erro vinculadas ao texto fonte, onde estas mensagens são produzidas em fases posteriores à análise léxica.

TIPO- é um valor numérico que determina a classe do lexema encontrado. Este valor é atribuído pelo transdutor finito reconhecedor. Os valores possíveis para TIPO são determinados durante a definição do filtro léxico. O TIPO será utilizado em fases posteriores à análise léxica para distinguir lexemas diferentes possuindo símbolos iguais. Além disso, permite tomar-se decisões durante a análise sintática sem que se tenha que recorrer ao símbolo contido na tabela de símbolos.

ID- é um valor numérico em correspondência um para um com o símbolo correspondente ao lexema reconhecido. O valor de ID é determinado pela tabela de símbolos.

Além desses resultados produzidos a cada ativação do filtro léxico efetuada pelo usuário, é produzido ainda um arquivo contendo a tabela de símbolos. Desta forma, poder-se-á ter acesso aos símbolos contidos na tabela, através de chamadas a rotinas específicas, tanto durante a filtragem léxica, quanto em fases posteriores.

Lexemas de erro têm sempre TIPO=0 e ID é a identificação de uma mensagem de erro explicando o erro havido. Os lexemas de erro são gerados imediatamente ao encontrar-se a condição de erro. Em muitos casos a condição de erro é recuperável em outros será final. No caso de uma condição de erro recuperável, a unidade léxica sendo reconhecida poderá continuar a ser reconhecida. Desta forma é possível continuar-se com a análise léxica, sintática e semântica, sem que o erro detectado provoque a geração de outros a nível sintático ou semântico. É claro que a condição de erro a nível léxico não é corrigida, no sentido estrito da palavra, porém é "remendada" o suficiente para que se possa prosseguir sem uma geração de erros consequência do erro encontrado.

Através da consulta ao campo ID e/ou aos caracteres iniciais da mensagem de erro, o usuário poderá determinar se se trata de um erro recuperável ou final. Em se tratando de um erro recuperável, o usuário poderá passar controle ao ponto de entrada de recuperação de erro (\$RECUPERE_ERRO()). Desta forma o reconhecimento da unidade léxica onde foi encontrado o erro é continuado. Caso o filtro léxico seja ativado através de seu ponto de entrada normal (\$FILTRO_LEXICO()), é iniciado o reconhecimento de uma nova unidade léxica, sendo perdida aquela em que ocorreu o erro. Caso o usuário ative o filtro léxico através do ponto de entrada de recuperação de erro quando o último lexema encontrado na instância ativa não havia sido lexema de erro recuperável, será gerado um lexema de erro final. Isto pode ocorrer se o último lexema encontrado nesta instância do filtro léxico tenha sido um lexema normal ou de erro final. Cabe frisar que é de inteira responsabilidade do usuário identificar se o lexema retornado é normal, de erro recuperável, ou de erro final. É seguinte a convenção adotada pelo gerador de filtros léxicos

- a- erros recuperáveis têm como primeiro caracter da mensagem de erro um caracter diferente de asterisco (*).
- b- erros finais têm como primeiro caracter da mensagem de erro o caracter asterisco (*).

As mensagens de erro poderão ser geradas devido a inconsistências reconhecidas pelo usuário ou pelo próprio filtro léxico. As mensagens de erro do usuário são especificadas durante a definição do filtro léxico (\$MENSAGEM(...) e \$ABORTE(...)). As mensagens de erro do sistema serão descritas em outra publicação (Manual de Especificação Física) sendo redigida no momento.

Além dos lexemas, é produzido, também, um arquivo contendo a tabela de símbolos. Este arquivo é de acesso direto, sendo que seu formato será descrito em outra publicação (Manual de Especificação Física) sendo redigida no momento. Ao encontrar o caracter especial \$FIM, o arquivo da tabela de símbolos é atualizado automaticamente, garantindo-se que esteja completo após esta atualização.

Em fases de processamento posteriores à filtragem léxica, este arquivo poderá ser acessado através de um conjunto de rotinas de acesso à tabela de símbolos. Este conjunto também é gerado pelo gerador de filtros léxicos.

2.4 Ativação.

O filtro léxico é ativado através de chamada a um de seus pontos de entrada. Os nomes destes pontos de entrada são definidos na seção <comunicação> da especificação do filtro léxico.

O ponto de entrada normal (`$FILTRO_LEXICO()`) sempre inicia o reconhecimento de um novo lexema, mesmo que o último lexema retornado pela instância ativa do filtro léxico tenha sido um lexema de erro recuperável.

O ponto de entrada de recuperação de erro (`$RECUPERE_ERRO()`) permite continuar no reconhecimento da unidade léxica em que foi encontrado um erro recuperável. Caso o último lexema retornado pela instância ativa do filtro léxico não seja lexema de erro recuperável, será gerado um lexema de erro final.

O retorno ocorrerá somente após o reconhecimento de um lexema. Como lexemas correspondem somente a unidades léxicas úteis, cabe ao usuário definir o filtro léxico de tal forma que não sejam induzidos "loops" infinitos de reconhecimento de unidades sintáticas inúteis. Atenção especial deve ser dada aos caracteres especiais `$FARQ`, `$FBUF`, `$FIM` e `$COND`, uma vez que estes podem facilmente induzir "loops" infinitos por poderem ser reconhecidos em sequências infinitas.

O reconhecimento do caracter `$ERRO` corresponde, sempre, a um lexema de erro. O usuário deverá tomar o cuidado de eliminar a causa do erro antes de reativar o filtro léxico, pois, de outra forma, poderá ser induzido um ciclo infinito de reconhecimento de condições de erro.

3 Especificação das Rotinas Componentes do Filtro Lógico.

Por razões de espaço disponível para este texto, apresentamos aqui somente um resumo da especificação do filtro lógico. Uma especificação completa encontra-se em [Staa78].

Cada filtro lógico específico é composto pelas rotinas:

- i- analisador lógico;
- ii- tabela de símbolos;
- iii- rotina de inicialização.

O analisador lógico tem por objetivos:

- i- decompor o texto fonte em unidades lógicas;
- ii- comunicar as porções úteis das cadeias de caracteres que perfazem unidades lógicas úteis para a área de comunicação da tabela de símbolos;
- iii- efetuar a comunicação com o exterior através da ativação de rotinas fornecidas pelo usuário, bem como através da ativação do próprio filtro lógico comandada pelo usuário.

A tabela de símbolos tem por objetivos:

- i- determinar um valor em correspondência um para um com os símbolos correspondentes a unidades lógicas úteis (lexemas);
- ii- permitir a salvaguarda dos símbolos contidos na tabela para acesso tanto durante a filtragem lógica, quanto durante fases posteriores;
- iii- prover rotinas de acesso a símbolos a partir de sua identificação.

A rotina de inicialização tem por objetivo criar instâncias de filtros lógicos e inicializar todas as variáveis a serem mantidas por uma determinada instância.

A análise lógica é a decomposição do fluxo de entrada em unidades lógicas.

A análise lógica é efetuada por um transdutor finito [Gries 1971, Hopcroft 1969, Aho 1972]. É função do transdutor finito:

- i- o reconhecimento de unidades lógicas, úteis ou não;

- ii- a eliminação de unidades léxicas inúteis;
- iii- a substituição de caracteres de partes do símbolo correspondentes a unidades léxicas úteis;
- iv- o reconhecimento de palavras reservadas e/ou pré-definidas;
- v- a ativação de rotinas fornecidas pelo usuário sempre que for efetuada uma transição requerendo esta ativação;
- vi- a intercomunicação com a tabela de símbolos;
- vii- a geração de lexemas indicativos de erro.

Símbolos correspondem à transdução de cadeias de caracteres das unidades léxicas encontradas. Um símbolo poderá ser qualquer cadeia de zero ou mais caracteres, sendo o seu tamanho virtualmente ilimitado. É função da tabela de símbolos estabelecer uma correspondência um para um entre a identificação do símbolo, campo ID do lexema, e o símbolo propriamente dito. Esta identificação permite recuperar-se o símbolo sem maiores dificuldades.

A tabela de símbolos é constituída por dois mapeamentos. Um mapeando cadeias de caracteres, os símbolos, sobre números naturais, e o outro mapeando números naturais sobre os símbolos correspondentes, ou, então, sobre o símbolo indefinido \emptyset . Formalmente:

$TS_1 : \text{Símbolos} \rightarrow N$
 $TS_2 : N \rightarrow \text{Símbolos} \cup \{ \emptyset \}$

A cada símbolo corresponde exatamente um número natural e a cada número natural corresponde exatamente um símbolo ou, então, o símbolo indefinido \emptyset , tal que:

$TS_2(TS_1(S)) = S$
se $TS_2(n) \neq \emptyset$ então $TS_1(TS_2(n)) = n$

Convém enfatizar que a finalidade única e exclusiva da tabela de símbolos é a de converter cadeias de caracteres em números naturais e vice-versa. Não é função da tabela de símbolos manter indicações quanto aos atributos (tipos) dos símbolos nela contidos. A tabela de atributos deve ser criada externamente ao filtro léxico [DeRemer74]. Tampouco é função do filtro léxico converter constantes numéricas em suas representações binárias internas. Caso isto seja desejado, o usuário deverá prover uma rotina para perfazer esta conversão.

O mapeamento TS_1 é chamado de instalação de símbolo, e é composto pelas seguintes operações:

- a- pesquisar a tabela de símbolos à procura do símbolo s a instalar;
- b- caso o símbolo a instalar já exista, devolver a identificação com que existe na tabela de símbolos;
- c- caso o símbolo a instalar ainda não exista, gerar uma identificação para ele, incluir a identificação e o símbolo na tabela de símbolos, e devolver a identificação criada.

O mapeamento TS_2 é chamado de recuperação de símbolo e é composto pelas seguintes operações:

- a- verificar se a identificação fornecida corresponde a um símbolo existente;
- b- caso corresponda, devolver o símbolo correspondente;
- c- caso não corresponda, devolver o símbolo indefinido.

4 Linguagem de Especificação de Filtros Léxicos.

Por limitações de espaço para este texto, apresentamos aqui somente uma fração da especificação da linguagem de especificação de filtros léxicos. A porção escolhida é bastante representativa do potencial desta linguagem. No apêndice é apresentada uma especificação completa de um filtro léxico escrito na linguagem de especificação. A descrição completa da linguagem de especificação de filtros léxicos encontra-se em [Staa78].

A linguagem de especificação de filtros léxicos está baseada na notação utilizada no compilador de compiladores COM-COM [Staa70]. Esta notação, por sua vez, é inspirada na notação de Kleene para linguagens regulares [Hopcroft69] e na notação Backus Naur Form (BNF) [Naur63].

Para a compreensão deste capítulo é assumido o leitor ter conhecimento de linguagens regulares, bem como das notações mencionadas. Detalhes quanto a estas notações poderão ser encontrados em [Hopcroft69, Gries71, Aho72, Bauer74].

A linguagem de especificação de filtros léxicos descreve de uma só vez tanto a linguagem regular reconhecida pelo transdutor (o conjunto de unidades léxicas permitidas), quanto as operações de transdução a serem efetuadas por este transdutor para produzir símbolos.

A linguagem de especificação de filtros l \acute{e} xicos \acute{e} definida em termos de opera \tilde{c} oes de reconhecimento, opera \tilde{c} oes de transdu \tilde{c} ao e de a \tilde{c} oes a serem efetuadas caso o reconhecimento chegue a um determinado ponto. O princ \acute{i} pio adotado \acute{e} semelhante ao adotado em SNOBOL4[Griswold71].

A leitura de \tilde{a} reas de texto de entrada, usualmente linhas, efetuada atrav \acute{e} s de rotina fornecida pelo usu \acute{a} rio, simula a exist \tilde{e} ncia de um fluxo de entrada cont \acute{i} nua e ininterrupto. Este efeito \acute{e} conseguido atrav \acute{e} s de chamadas \tilde{a} rotina de leitura do usu \acute{a} rio sempre que uma linha de texto de entrada exaurir. O fluxo ininterrupto \acute{e} conseguido atrav \acute{e} s da concatena \tilde{c} ao das linhas de texto de entrada simulando a exist \tilde{e} ncia entre elas do caracter especial \$FBUF. Sobre este fluxo de entrada \acute{e} movimentado um cursor. O reconhecimento de uma cadeia \acute{e} efetuado a partir do caracter sob este cursor. O reconhecimento somente suceder \tilde{a} , caso o fluxo de entrada contenha caracteres exatamente iguais aos especificados pela cadeia a reconhecer. As cadeias poder \tilde{a} o ser especificadas como sendo cadeias constantes, valores do tipo \$CADEIA contidos em vari \tilde{a} veis, valores do tipo \$CADEIA retornados por fun \tilde{c} oes do usu \acute{a} rio ou do sistema, ou, finalmente, um dos elementos de um conjunto de cadeias.

Dada a condi \tilde{c} ao de linguagem regular, o reconhecimento poder \tilde{a} ser efetuado de maneira determin \acute{i} stica, exceto quando a cadeia a reconhecer for resultado da avalia \tilde{c} ao de fun \tilde{c} oes ou conte \tilde{u} dos de vari \tilde{a} veis. Estas cadeias, por \tilde{e} m, n \tilde{a} o poder \tilde{a} o conter os caracteres especiais \$FIM, \$FBUF, \$FARQ, \$COND ou \$ERRO. Por esta raz \tilde{a} o, caso uma cadeia seja reconhec \acute{i} vel ela estar \tilde{a} integralmente contida no registro de entrada. Consequentemente, mesmo no caso de reconhecimento de cadeias vari \tilde{a} veis, este reconhecimento poder \tilde{a} ocorrer determin \acute{i} sticamente.

Cada especifica \tilde{c} ao de filtro l \acute{e} xico consta das seguintes quatro por \tilde{c} oes:

- i- cabe \tilde{c} alho - onde s \tilde{a} o dadas informa \tilde{c} oes sobre o t \acute{i} tulo, autor e datas, e descri \tilde{c} ao do filtro l \acute{e} xico;
- ii- comunica \tilde{c} ao - onde \acute{e} definida a rela \tilde{c} ao do filtro l \acute{e} xico a ser gerado com outros programas. A se \tilde{c} ao de comunica \tilde{c} ao descreve os pontos de entrada do filtro l \acute{e} xico gerado in \tilde{u} teis para fases posteriores, as refer \tilde{e} ncias a rotinas a serem providas pelo usu \acute{a} rio, os nomes de \tilde{a} reas externas e os nomes das se \tilde{c} oes de controle criadas pelo gerador de filtros l \acute{e} xicos;
- iii- defini \tilde{c} ao - onde s \tilde{a} o definidos os lexemas, os nomes de unidades sint \acute{a} ticas, conjuntos de cadeias de caracteres e vari \tilde{a} veis a serem utilizadas pelo filtro l \acute{e} xico;

- iv- transdutor - onde são definidas as unidades lêmicas reconhecíveis, bem como as transduções a serem efetuadas.

É seguinte a sintaxe para a definição de filtros lêmicos:

```
<filtro lêmico> ::= <cabeçalho>
                    <comunicação>
                    <definição>
                    <transdutor>
                    ($FIM | "$OUTRO" <filtro lêmico>).
```

Poderão ser definidos vários filtros lêmicos em uma somente aplicação. Estas definições deverão figurar em sequência, separadas pela palavra \$OUTRO. O término da última especificação de filtro lêmico é sinalizada pelo caracter especial \$FIM (ponteiro nulo retornado pela rotina de leitura).

4.2 Unidades Lêmicas da Linguagem de Especificação.

São seguintes os elementos lêmicos da linguagem de especificação de filtros lêmicos:

- i- cadeias - sequências de caracteres. Cadeias descrevem itens a reconhecer, bem como incorporações aos símbolos a serem efetuadas por transdução;
- ii- números - sequências de dígitos. Números são sempre maiores ou iguais a zero;
- iii- não terminais - elementos gramaticais que descrevem conjuntos de subsequências válidas que poderão ocorrer no interior de unidades lêmicas;
- iv- nomes - palavras reservadas ou palavras definidas pelo usuário, possuindo significado especial próprio;
- v- operadores e sinais de pontuação;
- vi- comentários e sequências de um ou mais caracteres em branco - texto fornecido pelo usuário para tornar a especificação mais legível e compreensível para o leitor. Comentários e sequências de caracteres em branco são ignorados pelo gerador de filtro lêmicos, servindo somente como separadores de elementos lêmicos. O uso de separador é mandatório somente se a sua ausência puder causar ambiguidades.

Descreveremos a seguir os elementos que constituem a linguagem de especificação de filtros lêmicos. A descrição utiliza a própria linguagem de especificação para definir os elementos lêmicos.

O texto a seguir descreve os lexemas utilizados pelo gerador de filtros l \acute{e} xicos. Estes lexemas est \tilde{a} o descritos na pr \acute{o} pria linguagem de especifica \tilde{c} o de filtros l \acute{e} xicos. Desta forma tem-se uma defini \tilde{c} o precisa, se bem que um pouco dif \acute{i} cil de compreender em alguns casos, pois necessita-se saber o que se est \tilde{a} tentando apreender...

No gerador de filtros l \acute{e} xicos as \acute{a} reas de texto de entrada correspondem a linhas de texto fonte. Chamaremos, ent \tilde{a} o, de linha as \acute{a} reas de texto fonte.

4.2.1 <Cadeia>.

Cadeias definem sequ \hat{e} ncias de zero ou mais caracteres que poder \tilde{a} o ser utilizadas para reconhecimento ou para a incorpora \tilde{c} o aos s \acute{i} mbois atrav \acute{e} s de transdu \tilde{c} o. O in \acute{i} cio das cadeias de caracteres \acute{e} sinalizado pelo caracter aspa dupla (""). O final da cadeia de caracteres \acute{e} sinalizado por outro caracter aspa dupla n \tilde{a} o precedido pelo caracter asterisco (*) ou por fim de linha de texto de entrada (\$FBUF). Cadeias s \tilde{a} o definidas da seguinte forma:

```
<cadeia> ::=
  ( <cadeia normal> |
    <caracter especial reconhec $\acute{i}$ vel> ).

<caracter especial reconhec $\acute{i}$ vel> ::=
  <nome> $NO_CONJUNTO (" $FIM" | " $FARQ" | " $FBUF" ).

<cadeia normal> ::= $ELIMINE "*" <corpo da cadeia>
  $ELIMINE <fim da cadeia> .

<fim da cadeia> ::=
  ( "*" |
    $FBUF $LE_LINHA() ).

<corpo da cadeia> ::= 0-(
  ( ~("*" | "~") |
    $ELIMINE "*" <opera $\tilde{c}$ o de cadeia> ) ).

/* <corpo da cadeia> termina somente se for encontrado
/* aspa dupla (") ou um caracter especial que  $\acute{e}$  $FBUF

<opera $\tilde{c}$ o de cadeia> ::=
  ( "*" |
    "*" |
    $TRADUZA $HEXA_BYTE(<byte hexadecimal>) |
    $ELIMINE <concatenador> |
    $MENSAGEM("0001Erro de opera $\tilde{c}$ o em cadeia.") ).

<byte hexadecimal> ::= 2 ($DIGITO_HEXADECIMAL).
```

```
<concatenador> ::= "-" $LE_LINHA() <concatenação>.  
<concatenação> ::= 0-(" "  
  ("*" |  
  $MENSAGEM("0002Erro de concatenação de cadeia."  
  $ENCONTROU <cadeia> ).
```

Durante o reconhecimento de uma cadeia, poderão ser encontrados comandos que atuam sobre esse reconhecimento. Através destes comandos torna-se possível a introdução de caracteres que não figurem no teclado do dispositivo de entrada utilizado para gerar a especificação; a introdução de caracteres com significado especial; e a geração de cadeias maiores do que as limitações impostas pelo equipamento de entrada utilizado.

Todos os comandos que atuam sobre o reconhecimento de cadeias são sinalizados pelo caracter de escape asterisco (*). São seguintes os comandos permitidos:

- *" introdução do caracter aspa dupla no corpo da cadeia;
- ** introdução do caracter asterisco no corpo da cadeia;
- *HH introdução no corpo da cadeia do caracter cuja representação em hexadecimal é HH;
- *- concatenação de linhas de texto de entrada. A linha corrente é terminada imediatamente ao encontrar-se este par de caracteres. A leitura da cadeia é reiniciada imediatamente após ter sido encontrado um caracter aspa dupla na linha de texto de entrada logo a seguir. Caso este caracter não seja encontrado, ou não seja o primeiro caracter não branco encontrado, é emitida uma mensagem de erro recuperável. Após a mensagem de erro é terminada a leitura da <cadeia>.

4.2.2 <Não Terminal>.

<Não terminal> é uma sequência de letras, dígitos, caracteres em branco e caracteres do conjunto :\$, @, #, _ : compreendidos entre os parênteses angulares < e >. O primeiro caracter do <não terminal> não deve ser dígito. Letras minúsculas serão traduzidas automaticamente para letras maiúsculas. Sequências de caracteres em branco serão traduzidas para exatamente um caracter em branco, exceto sequências logo após ao abre parêntese angular (<) e imediatamente antes do fecha parêntese angular (>) que serão eliminadas. O <não terminal> deve estar todo compreendido em uma linha de texto fonte.

"ABC"	a cadeia de corpo ABC
" "	a cadeia com corpo vazio
"*"	a cadeia de corpo "
"**"	a cadeia de corpo *
"*Fla***F2b*"*F3"	a cadeia de corpo la*2b"3
"AB*- "CD*- "EF"	a cadeia de corpo ABCDEF
"ABC" \$FBUF (Obs. \$ caracter em branco, \$FBUF fim de linha)	a cadeia de corpo ABC
"AB*- X"CD"	a mensagem de erro 0002

Figura 4.1 Exemplos de <cadeia>.

É seguinte a sintaxe dos <não terminais>:

CONJUNTO CARACTERES DE NOME = (\$LETRA | \$DIGITO |
"\$" | "@" | "#" | "_").

<não terminal> ::= \$ELIMINE "<"
(\$TRADUZA \$MAIUSCULAS(<corpo não terminal>)
\$ELIMINE <fim não terminal> |
\$MENSAGEM("0003Não terminal errado.")).

<corpo não terminal> ::= \$ELIMINE <brancos>
<resto corpo> .

<resto corpo> ::=
CARACTERES DE NOME ^(\$DIGITO)
0-((CARACTERES DE NOME |
" " \$ELIMINE <brancos>)).

<fim não terminal> ::=
(\$SE(\$ULTIMO TRANSFERIDO = " ") \$RETIRE(1) |
" ") ">" .

<brancos> ::= 0-(" ") .

```

<não terminal>
< não terminal >
< NãO TeRmInAl >

```

correspondem todos ao não terminal de corpo:
"NÃO TERMINAL"

Figura 4.2 Exemplos de <não terminal>.

4.2.3 <Nome>.

Um <nome> é uma sequência de letras, dígitos e caracteres do conjunto : \$, @, #, _:. O primeiro caracter de um <nome> não deve ser dígito. Letras minúsculas serão convertidas automaticamente para letras maiúsculas. Não podem existir caracteres em branco nos nomes. Nomes serão utilizados para identificar variáveis, funções, operações, conjuntos e comunicadores.

Os nomes pré-definidos pelo gerador de filtros léxicos e que principiam pelo caracter \$, são reservados. Nomes utilizados pelo gerador de filtros léxicos que não principiem pelo caracter \$ poderão ser redefinidos pelo usuário a qualquer momento. O usuário poderá definir nomes principiando pelo caracter \$, desde que não sejam iguais a nomes pré-definidos pelo gerador de filtros léxicos. Isto, porém, não é boa prática, uma vez que poderá confundir o leitor da especificação do filtro léxico.

É seguinte a sintaxe para <nomes>:

```

CONJUNTO CARACTERES DE NOME = ( $LETRA | $DIGITO |
"$" | "@" | "#" | "_" ).

```

```

<nome> ::= $TRADUZA $MAIUSCULAS(<corpo do nome>).

```

```

<corpo do nome> ::= CARACTERES DE NOME ^($DIGITO)
0-( CARACTERES _DÉ _NÔME ).

```

```

Ab_C           equivale a: AB_C
$#@_X
$
$LETRA
$letra        equivale a: $LETRA

```

Figura 4.3 Exemplos de <nomes>.

4.2.4 <Comentário> e <Sequência em Branco>.

<Comentário> é qualquer sequência de caracteres iniciada pelo par de caracteres /*, e terminada ou pelo par */ ou por fim de linha encontrado no fluxo de entrada.

<Sequência em branco> é uma sequência de um ou mais caracteres em branco e/ou fim de linha (\$FBUF).

<Comentário> e <sequência em branco> poderão ocorrer em qualquer lugar entre elementos léxicos da linguagem de especificação de filtros léxicos. Seu único efeito durante a filtragem léxica de uma especificação de filtro léxico é de delimitar componentes léxicos desta especificação, de resto são ignorados pelo gerador de filtros léxicos.

Cada <comentário> está compreendido no máximo em uma linha. Por outro lado, uma <sequência em branco> poderá ocupar seções de uma linha, linhas inteiras, podendo, inclusive, continuar na próxima linha.

É seguinte a sintaxe de <comentário> e <sequência em branco>:

```
<comentário> ::= "/*" 0-( ( ~("**") | "**" ~("/") )  
                ( "**" | "" ) ("/" | $FBUF $LE_LINHA() ) ) .
```

```
<sequência em branco> ::= 1-( " " | $FBUF $LE_LINHA() ) .
```

4.2.5 <Número> e Operadores.

Números são sequências de um ou mais dígitos decimais. Números são sempre maiores ou iguais a zero. É seguinte a sintaxe de <número>:

```
<número> ::= 1-($DIGITO).
```

O conjunto de operadores será definido à medida que forem apresentadas as construções permitidas na linguagem de especificação de filtros léxicos. No apêndice A é apresentada a especificação completa do filtro léxico utilizado pelo gerador de filtros léxicos. Neste apêndice são apresentados, também, os operadores existentes na linguagem de especificação de filtros léxicos.

5 Epílogo.

Uma especificação dificilmente se presta para uma seção de conclusões, em particular se a especificação é truncada tal como a presente. Apresentaremos então alguns argumentos justificando o esforço despendido no projeto e implementação.

Na PUC/RJ são diversas as linguagens sendo implementadas e/ou sendo projetadas para implementação no futuro próximo. Só isto já bastaria para justificar o desenvolvimento de uma ferramenta de suporte à implementação de processadores de linguagem.

Sabe-se hoje, também, que em muitos casos sistemas são fracassos devido a uma modalidade de comunicação homem-máquina inadequada. O gerador de filtros léxicos poderá servir, então, para o suporte de um laboratório experimental onde desenvolver novas técnicas de comunicação homem-máquina. Projetos dessa natureza também estão em andamento na PUC/RJ.

Referências Bibliográficas.

- Aho, A.V.; Ullman, J.D.
The Theory of Parsing, Translation and Compiling; Volume 1: Parsing, 1972; Volume 2: Compiling, 1973; Prentice Hall
- Bauer, F.L.; Eickel, J. editores
Compiler Construction, An Advanced Course; Springer Verlag, Lecture Notes in Computer Science no. 21, Berlin; 1974
- Brown, M.P.
Estudo Comparativo de Tabelas de Símbolos; Tese de Mestrado, Departamento de Informática, Pontifícia Universidade Católica, Rio de Janeiro; 1971
- DeRemer, F.L.
'Lexical Analysis'; em [Bauer74]; 1974; pags 109-120
- Gries, D:
Compiler Construction for Digital Computers; John Wiley and Sons, New York; 1971
- Griswold, R.E.; Poage, J.F.; Polonsky, I.P.
The SNOBOL4 Programming Language; Prentice Hall, Englewood Cliffs, New Jersey; Segunda Edição, 1971
- Hopcroft, J.E.; Ullman, J.D.
Formal Languages and Their Relation to Automata; Addison Wesley Publishing, Co., Reading, Massachusetts; 1969
- Johnson, W.R.; Porter, J.S.; Ackley, S.I.; Ross, D.T.
'Automatic Generation of Efficient Lexical Processors Using Finite State Techniques'; em Communications of the Association for Computing Machinery 11(12), New York; 1968; pags 805-813
- Knuth, D.E.
The Art of Computer Programming, Volume 3: Sorting and Searching; Addison Wesley Publishing Co., Reading, Massachusetts, 1973
- Naur, P. editor
'Revised Report on the Algorithmic Language ALGOL60'; em Numerische Mathematik 4; Springer Verlag, Berlin; 1963; pags 420-453
- Staa, A. v.
COMCOM - Compilador de Compiladores; Grupo de Aplicações, Rio Datacentro, Pontifícia Universidade Católica, Rio de Janeiro, 1970.

Staa, A. v.

Especificação do Projeto COMCOM2; Monografias em Ciência da Computação no. 16/77, Departamento de Informática, Pontifícia Universidade Católica, Rio de Janeiro; Outubro 1977.

Staa, A. v.

Gerador de Filtros Léxicos, Manual de Especificação Lógica; Monografias em Ciência da Computação, Departamento de Informática, Pontifícia Universidade Católica, Rio de Janeiro; 1978