```
**********************************************************
    8630 Fri Jan  3 08:20:33 2014
new/usr/src/cmd/sgs/lex/common/main.c
2926 lex ignores -Y
**********************************************************
   1 /*
   2  * CDDL HEADER START
   3  *
   4  * The contents of this file are subject to the terms of the
   5  * Common Development and Distribution License (the "License").
   6  * You may not use this file except in compliance with the License.
   7  *
   8  * You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE
   9  * or http://www.opensolaris.org/os/licensing.
  10  * See the License for the specific language governing permissions
  11  * and limitations under the License.
  12  *
  13  * When distributing Covered Code, include this CDDL HEADER in each
  14  * file and include the License file at usr/src/OPENSOLARIS.LICENSE.
  15  * If applicable, add the following below this CDDL HEADER, with the
  16  * fields enclosed by brackets "[]" replaced with your own identifying
  17  * information: Portions Copyright [yyyy] [name of copyright owner]
  18  *
  19  * CDDL HEADER END
  20  */
  21 /*
  22  * Copyright (c) 2014 Gary Mills
  23  *
  24  * Copyright 2008 Sun Microsystems, Inc.  All rights reserved.
  25  * Use is subject to license terms.
  26  */

  28 /* Copyright (c) 1988 AT&T */
  29 /* All Rights Reserved */

  31 /* Copyright 1976, Bell Telephone Laboratories, Inc. */

  31 #pragma ident   "%Z%%M% %I%     %E% SMI"

  33 #include <string.h>
  34 #include "once.h"
  35 #include "sgs.h"
  36 #include <locale.h>
  37 #include <limits.h>

  39 static wchar_t  L_INITIAL[] = {'I', 'N', 'I', 'T', 'I', 'A', 'L', 0};
  40 static void get1core(void);
  41 static void free1core(void);
  42 static void get2core(void);
  43 static void free2core(void);
  44 static void get3core(void);
  45 #ifdef DEBUG
  46 static void free3core(void);
  47 #endif

  49 int
  50 main(int argc, char **argv)
  51 {
  52         int i;
  53         int c;
  54         char *apath = NULL;
  55         char *ypath;
  54         char *path = NULL;
  56         Boolean eoption = 0, woption = 0;

  58         sargv = argv;
```

```
  59         sargc = argc;
  60         (void) setlocale(LC_ALL, "");
  61 #ifdef DEBUG
  62         while ((c = getopt(argc, argv, "dyctvnewVQ:Y:")) != EOF) {
  63 #else
  64         while ((c = getopt(argc, argv, "ctvnewVQ:Y:")) != EOF) {
  65 #endif
  66                 switch (c) {
  67 #ifdef DEBUG
  68                 case 'd':
  69                         debug++;
  70                         break;
  71                 case 'y':
  72                         yydebug = TRUE;
  73                         break;
  74 #endif
  75                 case 'V':
  76                         (void) fprintf(stderr, "lex: %s %s\n",
  77                             (const char *)SGU_PKG,
  78                             (const char *)SGU_REL);
  79                         break;
  80                 case 'Q':
  81                         v_stmp = optarg;
  82                         if (*v_stmp != 'y' && *v_stmp != 'n')
  83                                 error(
  84                                 "lex: -Q should be followed by [y/n]");
  85                         break;
  86                 case 'Y':
  87                         apath = (char *)malloc(strlen(optarg) +
  86                         path = (char *)malloc(strlen(optarg) +
  88                             sizeof ("/nceucform") + 1);
  89                         if (apath == NULL)
  90                                 error("No available memory "
  91                                     "for directory name.");
  92                         else
  93                                 apath = strcpy(apath, optarg);
  88                         path = strcpy(path, optarg);
  94                         break;
  95                 case 'c':
  96                         ratfor = FALSE;
  97                         break;
  98                 case 't':
  99                         fout = stdout;
 100                         break;
 101                 case 'v':
 102                         report = 1;
 103                         break;
 104                 case 'n':
 105                         report = 0;
 106                         break;
 107                 case 'w':
 108                 case 'W':
 109                         woption = 1;
 110                         handleeuc = 1;
 111                         widecio = 1;
 112                         break;
 113                 case 'e':
 114                 case 'E':
 115                         eoption = 1;
 116                         handleeuc = 1;
 117                         widecio = 0;
 118                         break;
 119                 default:
 120                         (void) fprintf(stderr,
 121                         "Usage: lex [-ewctvnV] [-Y directory] "
 122                         "[-Q(y/n)] [file]\n");
```

```
116                                 "Usage: lex [-ewctvnVY] [-Q(y/n)] [file]\n");
123                                 exit(1);
124                         }
125         }
126         if (woption && eoption) {
127                 error(
128                 "You may not specify both -w and -e simultaneously.");
129         }
130         no_input = argc - optind;
131         if (no_input) {
132                 /* XCU4: recognize "-" file operand for stdin */
133                 if (strcmp(argv[optind], "-") == 0)
134                         fin = stdin;
135                 else {
136                         fin = fopen(argv[optind], "r");
137                         if (fin == NULL)
138                                 error(
139                                 "Can't open input file -- %s", argv[optind]);
140                 }
141         } else
142                 fin = stdin;

144         /* may be gotten: def, subs, sname, schar, ccl, dchar */
145         (void) gch();

147         /* may be gotten: name, left, right, nullstr, parent */
148         get1core();

150         scopy(L_INITIAL, sp);
151         sname[0] = sp;
152         sp += slength(L_INITIAL) + 1;
153         sname[1] = 0;

155         /* XCU4: %x exclusive start */
156         exclusive[0] = 0;

158         if (!handleeuc) {
159                 /*
160                  * Set ZCH and ncg to their default values
161                  * as they may be needed to handle %t directive.
162                  */
163                 ZCH = ncg = NCH; /* ncg behaves as constant in this mode. */
164         }

166         /* may be disposed of: def, subs, dchar */
167         if (yyparse())
168                 exit(1);        /* error return code */

170         if (handleeuc) {
171                 ncg = ncgidtbl * 2;
172                 ZCH = ncg;
173                 if (ncg >= MAXNCG)
174                         error(
175                         "Too complex rules -- requires too many char groups.");
176                 sortcgidtbl();
177         }
178         repbycgid(); /* Call this even in ASCII compat. mode. */

180         /*
181          * maybe get:
182          *              tmpstat, foll, positions, gotof, nexts,
183          *              nchar, state, atable, sfall, cpackflg
184          */
185         free1core();
186         get2core();
187         ptail();
```

```
188         mkmatch();
189 #ifdef DEBUG
190         if (debug)
191                 pccl();
192 #endif
193         sect  = ENDSECTION;
194         if (tptr > 0)
195                 cfoll(tptr-1);
196 #ifdef DEBUG
197         if (debug)
198                 pfoll();
199 #endif
200         cgoto();
201 #ifdef DEBUG
202         if (debug) {
203                 (void) printf("Print %d states:\n", stnum + 1);
204                 for (i = 0; i <= stnum; i++)
205                         stprt(i);
206         }
207 #endif
208         /*
209          * may be disposed of:
210          *              positions, tmpstat, foll, state, name,
211          *              left, right, parent, ccl, schar, sname
212          * maybe get:   verify, advance, stoff
213          */
214         free2core();
215         get3core();
216         layout();
217         /*
218          * may be disposed of:
219          *              verify, advance, stoff, nexts, nchar,
220          *              gotof, atable, ccpackflg, sfall
221          */

223 #ifdef DEBUG
224         free3core();
225 #endif

227         if (handleeuc) {
228                 if (ratfor)
229                         error("Ratfor is not supported by -w or -e option.");
230                 ypath = EUCNAME;
224                 path = EUCNAME;
231         }
232         else
233                 ypath = ratfor ? RATNAME : CNAME;
227                 path = ratfor ? RATNAME : CNAME;

235         if (apath != NULL)
236                 ypath = strcat(apath, strrchr(ypath, '/'));
237         fother = fopen(ypath, "r");
229         fother = fopen(path, "r");
238         if (fother == NULL)
239                 error("Lex driver missing, file %s", ypath);
231                 error("Lex driver missing, file %s", path);
240         while ((i = getc(fother)) != EOF)
241                 (void) putc((char)i, fout);
242         (void) fclose(fother);
243         (void) fclose(fout);
244         free(apath);
245         if (report == 1)
246                 statistics();
247         (void) fclose(stdout);
248         (void) fclose(stderr);
249         return (0);     /* success return code */
```

```
 250 }
```
_____unchanged_portion_omitted_

     1 '\" te
     2 **.\" Copyright (c) 2014 Gary Mills**
     3 .\"  Copyright (c) 1992, X/Open Company Limited  All Rights Reserved  Portions C
     4 .\" Sun Microsystems, Inc. gratefully acknowledges The Open Group for permission
     5 .\" http://www.opengroup.org/bookstore/.
     6 .\" The Institute of Electrical and Electronics Engineers and The Open Group, ha
     7 .\"  This notice shall appear on any product containing this material.
     8 .\" The contents of this file are subject to the terms of the Common Development
     9 .\" You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE or http:
    10 .\" When distributing Covered Code, include this CDDL HEADER in each file and in
    11 **.TH LEX 1 "Jan 1, 2014"**
    10 *.TH LEX 1 "Aug 22, 1997"*
    12 .SH NAME
    13 lex \- generate programs for lexical tasks
    14 .SH SYNOPSIS
    15 .LP
    16 .nf
    17 **\fBlex\fR [\fB-cntv\fR] [\fB-e\fR | \fB-w\fR] [\fB-V\fR \fB-Q\fR [y | n]] [\fB-Y**
    16 *\fBlex\fR [\fB-cntv\fR] [\fB-e\fR | \fB-w\fR] [\fB-V\fR \fB-Q\fR [y | n]] [\fIfi*
    18 .fi

    20 .SH DESCRIPTION
    21 .sp
    22 .LP
    23 The \fBlex\fR utility generates C programs to be used in lexical processing of
    24 character input, and that can be used as an interface to \fByacc\fR. The C
    25 programs are generated from \fBlex\fR source code and conform to the ISO C
    26 standard. Usually, the \fBlex\fR utility writes the program it generates to the
    27 file \fBlex.yy.c\fR. The state of this file is unspecified if \fBlex\fR exits
    28 with a non-zero exit status. See \fBEXTENDED DESCRIPTION\fR for a complete
    29 description of the \fBlex\fR input language.
    30 .SH OPTIONS
    31 .sp
    32 .LP
    33 The following options are supported:
    34 .sp
    35 .ne 2
    36 .na
    37 \fB\fB-c\fR \fR
    38 .ad
    39 .RS 12n
    40 Indicates C-language action (default option).
    41 .RE

    43 .sp
    44 .ne 2
    45 .na
    46 \fB\fB-e\fR \fR
    47 .ad
    48 .RS 12n
    49 Generates a program that can handle \fBEUC\fR characters (cannot be used with
    50 the \fB-w\fR option). \fByytext[\|]\fR is of type \fBunsigned char[\|]\fR.
    51 .RE

    53 .sp
    54 .ne 2
    55 .na
    56 \fB\fB-n\fR \fR
    57 .ad
    58 .RS 12n
    59 Suppresses the summary of statistics usually written with the \fB-v\fR option.

    60 If no table sizes are specified in the \fBlex\fR source code and the \fB-v\fR
    61 option is not specified, then \fB-n\fR is implied.
    62 .RE

    64 .sp
    65 .ne 2
    66 .na
    67 \fB\fB-t\fR \fR
    68 .ad
    69 .RS 12n
    70 Writes the resulting program to standard output instead of \fBlex.yy.c\fR.
    71 .RE

    73 .sp
    74 .ne 2
    75 .na
    76 \fB\fB-v\fR \fR
    77 .ad
    78 .RS 12n
    79 Writes a summary of \fBlex\fR statistics to the standard error. (See the
    80 discussion of \fBlex\fR table sizes under the heading \fBDefinitions in
    81 lex\fR.) If table sizes are specified in the \fBlex\fR source code, and if the
    82 \fB-n\fR option is not specified, the \fB-v\fR option may be enabled.
    83 .RE

    85 .sp
    86 .ne 2
    87 .na
    88 \fB\fB-w\fR \fR
    89 .ad
    90 .RS 12n
    91 Generates a program that can handle \fBEUC\fR characters (cannot be used with
    92 the \fB-e\fR option). Unlike the \fB-e\fR option, \fByytext[\|]\fR is of type
    93 \fBwchar_t[\|]\fR.
    94 .RE

    96 .sp
    97 .ne 2
    98 .na
    99 \fB\fB-V\fR \fR
   100 .ad
   101 .RS 12n
   102 Prints out version information on standard error.
   103 .RE

   105 .sp
   106 .ne 2
   107 .na
   108 \fB\fB\fR\fB-Q\fR\fB[y|n]\fR \fR
   109 .ad
   110 .RS 12n
   111 Prints out version information to output file \fBlex.yy.c\fR by using
   112 \fB-Qy\fR. The \fB-Qn\fR option does not print out version information and is
   113 the default.
   114 .RE

   116 **.sp**
   117 **.ne 2**
   118 **.na**
   119 **\fB\fB\fR\fB-Y\fR \fBdirectory\fR \fR**
   120 **.ad**
   121 **.RS 12n**
   122 **Designates an alternate directory that contains the driver files**
   123 **used by \fBlex\fR.**
   124 **.RE**

```
126 .SH OPERANDS
127 .sp
128 .LP
129 The following operand is supported:
130 .sp
131 .ne 2
132 .na
133 \fB\fIfile\fR \fR
134 .ad
135 .RS 9n
136 A pathname of an input file. If more than one such \fIfile\fR is specified, all
137 files will be concatenated to produce a single \fBlex\fR program. If no
138 \fIfile\fR operands are specified, or if a \fIfile\fR operand is \fB\(mi\fR,
139 the standard input will be used.
140 .RE

142 .SH OUTPUT
143 .sp
144 .LP
145 The \fBlex\fR output files are described below.
146 .SS "Stdout"
147 .sp
148 .LP
149 If the \fB-t\fR option is specified, the text file of C source code output of
150 \fBlex\fR will be written to standard output.
151 .SS "Stderr"
152 .sp
153 .LP
154 If the \fB-t\fR option is specified informational, error and warning messages
155 concerning the contents of \fBlex\fR source code input will be written to the
156 standard error.
157 .sp
158 .LP
159 If the \fB-t\fR option is not specified:
160 .RS +4
161 .TP
162 1.
163 Informational error and warning messages concerning the contents of
164 \fBlex\fR source code input will be written to either the standard output or
165 standard error.
166 .RE
167 .RS +4
168 .TP
169 2.
170 If the \fB-v\fR option is specified and the \fB-n\fR option is not
171 specified, \fBlex\fR statistics will also be written to standard error. These
172 statistics may also be generated if table sizes are specified with a \fB%\fR
173 operator in the \fBDefinitions\fR \fBin\fR \fBlex\fR section (see \fBEXTENDED
174 DESCRIPTION\fR), as long as the \fB-n\fR option is not specified.
175 .RE
176 .SS "Output Files"
177 .sp
178 .LP
179 A text file containing C source code will be written to \fBlex.yy.c\fR, or to
180 the standard output if the \fB-t\fR option is present.
181 .SH EXTENDED DESCRIPTION
182 .sp
183 .LP
184 Each input file contains \fBlex\fR source code, which is a table of regular
185 expressions with corresponding actions in the form of C program fragments.
186 .sp
187 .LP
188 When \fBlex.yy.c\fR is compiled and linked with the \fBlex\fR library (using
189 the \fB\fR\fB-l\fR\fB l\fR operand with \fBc89\fR or \fBcc\fR), the resulting
190 program reads character input from the standard input and partitions it into
191 strings that match the given expressions.
```

```
192 .sp
193 .LP
194 When an expression is matched, these actions will occur:
195 .RS +4
196 .TP
197 .ie t \(bu
198 .el o
199 The input string that was matched is left in \fIyytext\fR as a null-terminated
200 string; \fIyytext\fR is either an external character array or a pointer to a
201 character string. As explained in \fBDefinitions\fR in lex\fR, the type can be
202 explicitly selected using the \fB%array\fR or \fB%pointer\fR declarations, but
203 the default is \fB%array\fR.
204 .RE
205 .RS +4
206 .TP
207 .ie t \(bu
208 .el o
209 The external \fBint\fR \fIyyleng\fR is set to the length of the matching
210 string.
211 .RE
212 .RS +4
213 .TP
214 .ie t \(bu
215 .el o
216 The expression's corresponding program fragment, or action, is executed.
217 .RE
218 .sp
219 .LP
220 During pattern matching, \fBlex\fR searches the set of patterns for the single
221 longest possible match. Among rules that match the same number of characters,
222 the rule given first will be chosen.
223 .sp
224 .LP
225 The general format of \fBlex\fR source is:
226 .sp
227 .in +2
228 .nf
229 \fIDefinitions\fR
230 %%
231 \fIRules\fR
232 %%
233 \fIUser Subroutines\fR
234 .fi
235 .in -2

237 .sp
238 .LP
239 The first \fB%%\fR is required to mark the beginning of the rules (regular
240 expressions and actions); the second \fB%%\fR is required only if user
241 subroutines follow.
242 .sp
243 .LP
244 Any line in the \fBDefinitions\fR \fBin\fR \fBlex\fR section beginning with a
245 blank character will be assumed to be a C program fragment and will be copied
246 to the external definition area of the \fBlex.yy.c\fR file. Similarly, anything
247 in the \fBDefinitions\fR \fBin\fR \fBlex\fR section included between delimiter
248 lines containing only \fB%{\fR and \fB%}\fR will also be copied unchanged to
249 the external definition area of the \fBlex.yy.c\fR file.
250 .sp
251 .LP
252 Any such input (beginning with a blank character or within \fB%{\fR and
253 \fB%}\fR delimiter lines) appearing at the beginning of the \fIRules\fR section
254 before any rules are specified will be written to \fBlex.yy.c\fR after the
255 declarations of variables for the \fByylex\fR function and before the first
256 line of code in \fByylex\fR. Thus, user variables local to \fByylex\fR can be
257 declared here, as well as application code to execute upon entry to
```

```
 258 \fByylex\fR.
 259 .sp
 260 .LP
 261 The action taken by \fBlex\fR when encountering any input beginning with a
 262 blank character or within \fB%{\fR and \fB%}\fR delimiter lines appearing in
 263 the \fIRules\fR section but coming after one or more rules is undefined. The
 264 presence of such input may result in an erroneous definition of the \fByylex\fR
 265 function.
 266 .SS "Definitions in lex"
 267 .sp
 268 .LP
 269 \fBDefinitions\fR \fBin\fR \fBlex\fR appear before the first \fB%%\fR
 270 delimiter. Any line in this section not contained between \fB%{\fR and \fB%}\fR
 271 lines and not beginning with a blank character is assumed to define a \fBlex\fR
 272 substitution string. The format of these lines is:
 273 .sp
 274 .in +2
 275 .nf
 276 \fIname    substitute\fR
 277 .fi
 278 .in -2
 279 .sp
 281 .sp
 282 .LP
 283 If a \fIname\fR does not meet the requirements for identifiers in the ISO C
 284 standard, the result is undefined. The string \fIsubstitute\fR will replace the
 285 string \fI{\fR \fIname\fR \fI}\fR when it is used in a rule. The \fIname\fR
 286 string is recognized in this context only when the braces are provided and when
 287 it does not appear within a bracket expression or within double-quotes.
 288 .sp
 289 .LP
 290 In the \fBDefinitions\fR \fBin\fR \fBlex\fR section, any line beginning with a
 291 \fB%\fR (percent sign) character and followed by an alphanumeric word beginning
 292 with either \fBs\fR or \fBS\fR defines a set of start conditions. Any line
 293 beginning with a \fB%\fR followed by a word beginning with either \fBx\fR or
 294 \fBX\fR defines a set of exclusive start conditions. When the generated scanner
 295 is in a \fBs\fR state, patterns with no state specified will be also active;
 296 in a \fBx\fR state, such patterns will not be active. The rest of the line,
 297 after the first word, is considered to be one or more blank-character-separated
 298 names of start conditions. Start condition names are constructed in the same
 299 way as definition names. Start conditions can be used to restrict the matching
 300 of regular expressions to one or more states as described in \fBRegular
 301 expressions in lex\fR.
 302 .sp
 303 .LP
 304 Implementations accept either of the following two mutually exclusive
 305 declarations in the \fBDefinitions\fR \fBin\fR \fBlex\fR section:
 306 .sp
 307 .ne 2
 308 .na
 309 \fB\fB%array\fR \fR
 310 .ad
 311 .RS 13n
 312 Declare the type of \fIyytext\fR to be a null-terminated character array.
 313 .RE
 315 .sp
 316 .ne 2
 317 .na
 318 \fB\fB%pointer\fR \fR
 319 .ad
 320 .RS 13n
 321 Declare the type of \fIyytext\fR to be a pointer to a null-terminated character
 322 string.
 323 .RE
```

```
 325 .sp
 326 .LP
 327 \fBNote:\fR When using the \fB%pointer\fR option, you may not also use the
 328 \fByyless\fR function to alter \fIyytext\fR.
 329 .sp
 330 .LP
 331 \fB%array\fR is the default. If \fB%array\fR is specified (or neither
 332 \fB%array\fR nor \fB%pointer\fR is specified), then the correct way to make an
 333 external reference to \fIyytext\fR is with a declaration of the form:
 334 .sp
 335 .LP
 336 \fBextern char\fR\fI yytext\fR\fB[\|]\fR
 337 .sp
 338 .LP
 339 If \fB%pointer\fR is specified, then the correct external reference is of the
 340 form:
 341 .sp
 342 .LP
 343 \fBextern char *\fR\fIyytext\fR\fB;\fR
 344 .sp
 345 .LP
 346 \fBlex\fR will accept declarations in the \fBDefinitions in lex\fR section for
 347 setting certain internal table sizes. The declarations are shown in the
 348 following table.
 349 .sp
 350 .LP
 351 \fBTable\fR \fBSize\fR \fBDeclaration\fR \fBin\fR \fBlex\fR
 352 .sp
 354 .sp
 355 .TS
 356 box;
 357 c c c
 358 l l l .
 359 \fBDeclaration\fR        \fBDescription\fR        \fBDefault\fR
 360 _
 361 \fB%p\fR\fIn\fR Number of positions     2500
 362 \fB%n\fR\fIn\fR Number of states        500
 363 \fB%a\fR\fI n\fR        Number of transitions   2000
 364 \fB%e\fR\fIn\fR Number of parse tree nodes      1000
 365 \fB%k\fR\fIn\fR Number of packed character classes       10000
 366 \fB%o\fR\fIn\fR Size of the output array        3000
 367 .TE
 369 .sp
 370 .LP
 371 Programs generated by \fBlex\fR need either the \fB-e\fR or \fB-w\fR option to
 372 handle input that contains \fBEUC\fR characters from supplementary codesets. If
 373 neither of these options is specified, \fByytext\fR is of the type
 374 \fBchar[\|]\fR, and the generated program can handle only \fBASCII\fR
 375 characters.
 376 .sp
 377 .LP
 378 When the \fB-e\fR option is used, \fByytext\fR is of the type \fBunsigned\fR
 379 \fBchar[\|]\fR and \fByyleng\fR gives the total number of \fIbytes\fR in the
 380 matched string. With this option, the macros \fBinput()\fR,
 381 \fBunput(\fIc\fR)\fR, and \fBoutput(\fIc\fR)\fR should do a byte-based
 382 \fBI/O\fR in the same way as with the regular \fBASCII\fR \fBlex\fR. Two more
 383 variables are available with the \fB-e\fR option, \fByywtext\fR and
 384 \fByywleng\fR, which behave the same as \fByytext\fR and \fByyleng\fR would
 385 under the \fB-w\fR option.
 386 .sp
 387 .LP
 388 When the \fB-w\fR option is used, \fByytext\fR is of the type \fBwchar_t[\|]\fR
 389 and \fByyleng\fR gives the total number of \fIcharacters\fR in the matched
```

```
 390 string.  If you supply your own \fBinput()\fR, \fBunput(\fIc\fR)\fR, or
 391 \fBoutput(\fR\fIc\fR\fB)\fR macros with this option, they must return or accept
 392 \fBEUC\fR characters in the form of wide character (\fBwchar_t\fR). This allows
 393 a different interface between your program and the lex internals, to expedite
 394 some programs.
 395 .SS "Rules in lex"
 396 .sp
 397 .LP
 398 The \fBRules\fR \fBin\fR \fBlex\fR source files are a table in which the left
 399 column contains regular expressions and the right column contains actions (C
 400 program fragments) to be executed when the expressions are recognized.
 401 .sp
 402 .in +2
 403 .nf
 404 \fIERE action\fR
 405 \fIERE action\fR
 406 \&...
 407 .fi
 408 .in -2

 410 .sp
 411 .LP
 412 The extended regular expression (ERE) portion of a row will be separated from
 413 \fIaction\fR by one or more blank characters. A regular expression containing
 414 blank characters is recognized under one of the following conditions:
 415 .RS +4
 416 .TP
 417 .ie t \(bu
 418 .el o
 419 The entire expression appears within double-quotes.
 420 .RE
 421 .RS +4
 422 .TP
 423 .ie t \(bu
 424 .el o
 425 The blank characters appear within double-quotes or square brackets.
 426 .RE
 427 .RS +4
 428 .TP
 429 .ie t \(bu
 430 .el o
 431 Each blank character is preceded by a backslash character.
 432 .RE
 433 .SS "User Subroutines in lex"
 434 .sp
 435 .LP
 436 Anything in the user subroutines section will be copied to \fBlex.yy.c\fR
 437 following \fByylex\fR.
 438 .SS "Regular Expressions     in lex"
 439 .sp
 440 .LP
 441 The \fBlex\fR utility supports the set of Extended Regular Expressions (EREs)
 442 described on \fBregex\fR(5) with the following additions and exceptions to the
 443 syntax:
 444 .sp
 445 .ne 2
 446 .na
 447 \fB\fB\|.\|.\|.\fR \fR
 448 .ad
 449 .RS 14n
 450 Any string enclosed in double-quotes will represent the characters within the
 451 double-quotes as themselves, except that backslash escapes (which appear in the
 452 following table) are recognized. Any backslash-escape sequence is terminated by
 453 the closing quote. For example, "\|\e\|01""1" represents a single string: the
 454 octal value 1 followed by the character 1.
 455 .RE
```

```
 457 .sp
 458 .LP
 459 \fI<\fR\fIstate\fR\fI>\fR\fIr\fR
 460 .sp
 461 .ne 2
 462 .na
 463 \fB<\fIstate1\fR, \fIstate2\fR, \|.\|.\|.\|>\fIr\fR\fR
 464 .ad
 465 .sp .6
 466 .RS 4n
 467 The regular expression \fIr\fR will be matched only when the program is in one
 468 of the start conditions indicated by \fIstate\fR, \fIstate1\fR, and so forth.
 469 For more information, see \fBActions in lex\fR. As an exception to the
 470 typographical conventions of the rest of this document, in this case
 471 <\fIstate\fR> does not represent a metavariable, but the literal angle-bracket
 472 characters surrounding a symbol. The start condition is recognized as such only
 473 at the beginning of a regular expression.
 474 .RE

 476 .sp
 477 .ne 2
 478 .na
 479 \fB\fIr\fR/\fIx\fR \fR
 480 .ad
 481 .sp .6
 482 .RS 4n
 483 The regular expression \fIr\fR will be matched only if it is followed by an
 484 occurrence of regular expression \fIx\fR. The token returned in \fIyytext\fR
 485 will only match \fIr\fR. If the trailing portion of \fIr\fR matches the
 486 beginning of \fIx\fR, the result is unspecified. The \fIr\fR expression cannot
 487 include further trailing context or the \fB$\fR (match-end-of-line) operator;
 488 \fIx\fR cannot include the \fB^\fR (match-beginning-of-line) operator, nor
 489 trailing context, nor the \fB$\fR operator. That is, only one occurrence of
 490 trailing context is allowed in a \fBlex\fR regular expression, and the \fB^\fR
 491 operator only can be used at the beginning of such an expression. A further
 492 restriction is that the trailing-context operator \fB/\fR (slash) cannot be
 493 grouped within parentheses.
 494 .RE

 496 .sp
 497 .ne 2
 498 .na
 499 \fB\fB{\fR\fIname\fR\fB}\fR \fR
 500 .ad
 501 .sp .6
 502 .RS 4n
 503 When \fIname\fR is one of the substitution symbols from the \fIDefinitions\fR
 504 section, the string, including the enclosing braces, will be replaced by the
 505 \fIsubstitute\fR value. The \fIsubstitute\fR value will be treated in the
 506 extended regular expression as if it were enclosed in parentheses. No
 507 substitution will occur if \fB{\fR\fIname\fR\fB}\fR occurs within a bracket
 508 expression or within double-quotes.
 509 .RE

 511 .sp
 512 .LP
 513 Within an \fBERE,\fR a backslash character (\fB\|\e\e\fR, \fB\e\|a\fR,
 514 \fB\e\|b\fR, \fB\e\|f\fR, \fB\e\|n\fR, \fB\e\|r\fR, \fB\e\|t\fR, \fB\e\|v\fR)
 515 is considered to begin an escape sequence. In addition, the escape sequences in
 516 the following table will be recognized.
 517 .sp
 518 .LP
 519 A literal newline character cannot occur within an \fBERE;\fR the escape
 520 sequence \fB\e\|n\fR can be used to represent a newline character. A newline
 521 character cannot be matched by a period operator.
```

```
 522 .sp
 523 .LP
 524 \fBEscape Sequences in lex\fR
 525 .sp

 527 .sp
 528 .TS
 529 box;
 530 c c c
 531 c c c .
 532 Escape Sequences in lex
 533 _
 534 Escape Sequence Description     Meaning
 535 _
 536 \e\fIdigits\fR  T{
 537 A backslash character followed by the longest sequence of one, two or three octa
 538 T}      T{
 539 The character whose encoding is represented by the one-, two- or three-digit oct
 540 T}
 541 _
 542 \e\fBx\fR\fIdigits\fR    T{
 543 A backslash character followed by the longest sequence of hexadecimal-digit char
 544 T}      T{
 545 The character whose encoding is represented by the hexadecimal integer.
 546 T}
 547 _
 548 \e\fIc\fR       T{
 549 A backslash character followed by any character not described in this table.  (\
 550 T}      The character c, unchanged.
 551 .TE

 553 .sp
 554 .LP
 555 The order of precedence given to extended regular expressions for \fBlex\fR is
 556 as shown in the following table, from high to low.
 557 .sp
 558 .ne 2
 559 .na
 560 \fB\fBNote\fR: \fR
 561 .ad
 562 .RS 10n
 563 The escaped characters entry is not meant to imply that these are operators,
 564 but they are included in the table to show their relationships to the true
 565 operators. The start condition, trailing context and anchoring notations have
 566 been omitted from the table because of the placement restrictions described in
 567 this section; they can only appear at the beginning or ending of an \fBERE.\fR
 568 .RE

 570 .sp

 572 .sp
 573 .TS
 574 box;
 575 c c
 576 l l .
 577 ERE Precedence in lex
 578 _
 579 \fIcollation-related bracket symbols\fR \fB[= =]  [: :]  [. .]\fR
 580 \fIescaped characters\fR        \fB\e<\fR\fIspecial character\fR>
 581 \fIbracket expression\fR        \fB[ ]\fR
 582 \fIquoting\fR   \fB".\|.\|.\|."\fR
 583 \fIgrouping\fR  \fB()\fR
 584 \fIdefinition\fR        \fB{\fR\fIname\fR}
 585 \fIsingle-character RE duplication\fR   \fB* + ?\fR
 586 \fIconcatenation\fR
 587 \fIinterval expression\fR       \fB{\fR\fIm\fR,\fIn\fR}
```

```
 588 \fIalternation\fR       \fB|\fR
 589 .TE

 591 .sp
 592 .LP
 593 The \fBERE\fR anchoring operators (\fB\|^\fR and \fB$\fR\|) do not appear in
 594 the table. With \fBlex\fR regular expressions, these operators are restricted
 595 in their use: the \fB^\fR operator can only be used at the beginning of an
 596 entire regular expression, and the \fB$\fR operator only at the end. The
 597 operators apply to the entire regular expression. Thus, for example, the
 598 pattern (\fB^abc\|(def$\fR) is undefined; it can instead be written as two
 599 separate rules, one with the regular expression \fB^abc\fR and one with
 600 \fBdef$\fR, which share a common action via the special \fB|\fR action (see
 601 below). If the pattern were written \fB^abc|def$\fR, it would match either of
 602 \fBabc\fR or \fBdef\fR on a line by itself.
 603 .sp
 604 .LP
 605 Unlike the general \fBERE\fR rules, embedded anchoring is not allowed by most
 606 historical \fBlex\fR implementations. An example of embedded anchoring would be
 607 for patterns such as (^)foo($) to match \fBfoo\fR when it exists as a complete
 608 word. This functionality can be obtained using existing \fBlex\fR features:
 609 .sp
 610 .in +2
 611 .nf
 612 ^foo/[ \e\|n]\|
 613 " foo"/[ \e\|n]     /* found foo as a separate word */
 614 .fi
 615 .in -2

 617 .sp
 618 .LP
 619 Notice also that \fB$\fR is a form of trailing context (it is equivalent to
 620 \fB/\e\|n\fR and as such cannot be used with regular expressions containing
 621 another instance of the operator (see the preceding discussion of trailing
 622 context).
 623 .sp
 624 .LP
 625 The additional regular expressions trailing-context operator \fB/\fR (slash)
 626 can be used as an ordinary character if presented within double-quotes,
 627 \fB"\|/\|"\fR; preceded by a backslash, \fB\e\|/\fR; or within a bracket
 628 expression, \fB[\|/\|]\fR. The start-condition \fB<\fR and \fB>\fR operators
 629 are special only in a start condition at the beginning of a regular expression;
 630 elsewhere in the regular expression they are treated as ordinary characters.
 631 .sp
 632 .LP
 633 The following examples clarify the differences between \fBlex\fR regular
 634 expressions and regular expressions appearing elsewhere in this document. For
 635 regular expressions of the form \fIr\fR/\fIx\fR, the string matching \fIr\fR is
 636 always returned; confusion may arise when the beginning of \fIx\fR matches the
 637 trailing portion of \fIr\fR. For example, given the regular expression a*b/cc
 638 and the input \fBaaabcc\fR, \fIyytext\fR would contain the string \fBaaab\fR on
 639 this match. But given the regular expression x*/xy and the input \fBxxxy\fR,
 640 the token \fBxxx\fR, not \fBxx\fR, is returned by some implementations because
 641 \fBxxx\fR matches x*.
 642 .sp
 643 .LP
 644 In the rule ab*/bc, the b* at the end of \fIr\fR will extend \fIr\fR's match
 645 into the beginning of the trailing context, so the result is unspecified. If
 646 this rule were ab/bc, however, the rule matches the text \fBab\fR when it is
 647 followed by the text \fBbc\fR. In this latter case, the matching of \fIr\fR
 648 cannot extend into the beginning of \fIx\fR, so the result is specified.
 649 .SS "Actions in lex"
 650 .sp
 651 .LP
 652 The action to be taken when an \fBERE\fR is matched can be a C program fragment
 653 or the special actions described below; the program fragment can contain one or
```

654 more C statements, and can also include special actions. The empty C statement
655 \fB;\fR is a valid action; any string in the \fBlex.yy.c\fR input that matches
656 the pattern portion of such a rule is effectively ignored or skipped. However,
657 the absence of an action is not valid, and the action \fBlex\fR takes in such a
658 condition is undefined.
659 .sp
660 .LP
661 The specification for an action, including C statements and special actions,
662 can extend across several lines if enclosed in braces:
663 .sp
664 .in +2
665 .nf
666 ERE <one or more blanks> { program statement
667 program statement }
668 .fi
669 .in -2
670 .sp

672 .sp
673 .LP
674 The default action when a string in the input to a \fBlex.yy.c\fR program is
675 not matched by any expression is to copy the string to the output. Because the
676 default behavior of a program generated by \fBlex\fR is to read the input and
677 copy it to the output, a minimal \fBlex\fR source program that has just
678 \fB%%\fR generates a C program that simply copies the input to the output
679 unchanged.
680 .sp
681 .LP
682 Four special actions are available:
683 .sp
684 .in +2
685 .nf
686 |       ECHO;       REJECT;       BEGIN
687 .fi
688 .in -2
689 .sp

691 .sp
692 .ne 2
693 .na
694 \fB|\fR
695 .ad
696 .RS 12n
697 The action | means that the action for the next rule is the action for this
698 rule. Unlike the other three actions, | cannot be enclosed in braces or be
699 semicolon-terminated. It must be specified alone, with no other actions.
700 .RE

702 .sp
703 .ne 2
704 .na
705 \fB\fBECHO;\fR \fR
706 .ad
707 .RS 12n
708 Writes the contents of the string \fIyytext\fR on the output.
709 .RE

711 .sp
712 .ne 2
713 .na
714 \fB\fBREJECT;\fR \fR
715 .ad
716 .RS 12n
717 Usually only a single expression is matched by a given string in the input.
718 \fBREJECT\fR means "continue to the next expression that matches the current
719 input," and causes whatever rule was the second choice after the current rule

720 to be executed for the same input. Thus, multiple rules can be matched and
721 executed for one input string or overlapping input strings. For example, given
722 the regular expressions \fBxyz\fR and \fBxy\fR and the input \fBxyz\fR, usually
723 only the regular expression \fBxyz\fR would match. The next attempted match
724 would start after z. If the last action in the \fBxyz\fR rule is \fBREJECT\fR ,
725 both this rule and the \fBxy\fR rule would be executed. The \fBREJECT\fR action
726 may be implemented in such a fashion that flow of control does not continue
727 after it, as if it were equivalent to a \fBgoto\fR to another part of
728 \fByylex\fR. The use of \fBREJECT\fR may result in somewhat larger and slower
729 scanners.
730 .RE

732 .sp
733 .ne 2
734 .na
735 \fB\fBBEGIN\fR \fR
736 .ad
737 .RS 12n
738 The action:
739 .sp
740 \fBBEGIN\fR \fInewstate\fR\fB;\fR
741 .sp
742 switches the state (start condition) to \fInewstate\fR. If the string
743 \fInewstate\fR has not been declared previously as a start condition in the
744 \fBDefinitions\fR \fBin\fR \fBlex\fR section, the results are unspecified. The
745 initial state is indicated by the digit \fB0\fR or the token \fBINITIAL\fR.
746 .RE

748 .sp
749 .LP
750 The functions or macros described below are accessible to user code included in
751 the \fBlex\fR input. It is unspecified whether they appear in the C code output
752 of \fBlex\fR, or are accessible only through the \fB\fR\fB-l\fR\fB l\fR operand
753 to \fBc89\fR or \fBcc\fR (the \fBlex\fR library).
754 .sp
755 .ne 2
756 .na
757 \fB\fBint\fR \fByylex(void)\fR \fR
758 .ad
759 .RS 21n
760 Performs lexical analysis on the input; this is the primary function generated
761 by the \fBlex\fR utility. The function returns zero when the end of input is
762 reached; otherwise it returns non-zero values (tokens) determined by the
763 actions that are selected.
764 .RE

766 .sp
767 .ne 2
768 .na
769 \fB\fBint\fR \fByymore(void)\fR \fR
770 .ad
771 .RS 21n
772 When called, indicates that when the next input string is recognized, it is to
773 be appended to the current value of \fIyytext\fR rather than replacing it; the
774 value in \fIyyleng\fR is adjusted accordingly.
775 .RE

777 .sp
778 .ne 2
779 .na
780 \fB\fBint\fR\fIyyless(int\fR\fB n\fR\fI)\fR \fR
781 .ad
782 .RS 21n
783 Retains \fIn\fR initial characters in \fIyytext\fR, NUL-terminated, and treats
784 the remaining characters as if they had not been read; the value in
785 \fIyyleng\fR is adjusted accordingly.

```
 786 .RE

 788 .sp
 789 .ne 2
 790 .na
 791 \fB\fIint\fR \fBinput(void)\fR \fR
 792 .ad
 793 .RS 21n
 794 Returns the next character from the input, or zero on end-of-file. It obtains
 795 input from the stream pointer \fIyyin\fR, although possibly via an intermediate
 796 buffer. Thus, once scanning has begun, the effect of altering the value of
 797 \fIyyin\fR is undefined. The character read is removed from the input stream of
 798 the scanner without any processing by the scanner.
 799 .RE

 801 .sp
 802 .ne 2
 803 .na
 804 \fB\fIint\fR \fBunput(int\fR \fB\fIc\fR\fR\fB)\fR \fR
 805 .ad
 806 .RS 21n
 807 Returns the character \fIc\fR to the input; \fIyytext\fR and \fIyyleng\fR are
 808 undefined until the next expression is matched. The result of using \fIunput\fR
 809 for more characters than have been input is unspecified.
 810 .RE

 812 .sp
 813 .LP
 814 The following functions appear only in the \fBlex\fR library accessible through
 815 the \fB\fR\fB-l\fR\fB l\fR operand; they can therefore be redefined by a
 816 portable application:
 817 .sp
 818 .ne 2
 819 .na
 820 \fB\fIint\fR \fByywrap(void)\fR \fR
 821 .ad
 822 .sp .6
 823 .RS 4n
 824 Called by \fByylex\fR at end-of-file; the default \fByywrap\fR always will
 825 return 1. If the application requires \fByylex\fR to continue processing with
 826 another source of input, then the application can include a function
 827 \fByywrap\fR, which associates another file with the external variable
 828 \fBFILE\fR *\fIyyin\fR and will return a value of zero.
 829 .RE

 831 .sp
 832 .ne 2
 833 .na
 834 \fB\fIint\fR \fBmain(int\fR \fB\fIargc\fR,\fR \fBchar\fR \fB*\fIargv\fR\fR[\|])\fR
 835 \fR
 836 .ad
 837 .sp .6
 838 .RS 4n
 839 Calls \fByylex\fR to perform lexical analysis, then exits. The user code can
 840 contain \fBmain\fR to perform application-specific operations, calling
 841 \fByylex\fR as applicable.
 842 .RE

 844 .sp
 845 .LP
 846 The reason for breaking these functions into two lists is that only those
 847 functions in \fBlibl.a\fR can be reliably redefined by a portable application.
 848 .sp
 849 .LP
 850 Except for \fBinput\fR, \fBunput\fR and \fBmain\fR, all external and static
 851 names generated by \fBlex\fR begin with the prefix \fByy\fR or \fBYY\fR.
```

```
 852 .SH USAGE
 853 .sp
 854 .LP
 855 Portable applications are warned that in the \fBRules in lex\fR section, an
 856 \fBERE\fR without an action is not acceptable, but need not be detected as
 857 erroneous by \fBlex\fR. This may result in compilation or run-time errors.
 858 .sp
 859 .LP
 860 The purpose of \fBinput\fR is to take characters off the input stream and
 861 discard them as far as the lexical analysis is concerned. A common use is to
 862 discard the body of a comment once the beginning of a comment is recognized.
 863 .sp
 864 .LP
 865 The \fBlex\fR utility is not fully internationalized in its treatment of
 866 regular expressions in the \fBlex\fR source code or generated lexical analyzer.
 867 It would seem desirable to have the lexical analyzer interpret the regular
 868 expressions given in the \fBlex\fR source according to the environment
 869 specified when the lexical analyzer is executed, but this is not possible with
 870 the current \fBlex\fR technology. Furthermore, the very nature of the lexical
 871 analyzers produced by \fBlex\fR must be closely tied to the lexical
 872 requirements of the input language being described, which will frequently be
 873 locale-specific anyway. (For example, writing an analyzer that is used for
 874 French text will not automatically be useful for processing other languages.)
 875 .SH EXAMPLES
 876 .LP
 877 \fBExample 1 \fRUsing lex
 878 .sp
 879 .LP
 880 The following is an example of a \fBlex\fR program that implements a
 881 rudimentary scanner for a Pascal-like syntax:

 883 .sp
 884 .in +2
 885 .nf
 886 %{
 887 /* need this for the call to atof() below */
 888 #include <math.h>
 889 /* need this for printf(), fopen() and stdin below */
 890 #include <stdio.h>
 891 %}

 893 DIGIT    [0-9]
 894 ID       [a-z][a-z0-9]*
 895 %%

 897 {DIGIT}+        {
 898                         printf("An integer: %s (%d)\en", yytext,
 899                         atoi(yytext));
 900                         }

 902 {DIGIT}+"."{DIGIT}*     {
 903                         printf("A float: %s (%g)\en", yytext,
 904                         atof(yytext));
 905                         }

 907 if|then|begin|end|procedure|function        {
 908                         printf("A keyword: %s\en", yytext);
 909                         }

 911 {ID}            printf("An identifier: %s\en", yytext);

 913 "+"|"-"|"*"|"/"        printf("An operator: %s\en", yytext);

 915 "{"[^}\en]*"}"         /* eat up one-line comments */

 917 [ \et\en]+            /* eat up white space */
```

```
 919 \&.                          printf("Unrecognized character: %s\en", yytext);

 921 %%

 923 int main(int argc, char *argv[\|])
 924 {
 925                     ++argv, --argc;  /* skip over program name */
 926                     if (argc > 0)
 927                             yyin = fopen(argv[0], "r");
 928                     else
 929                     yyin = stdin;
 930
 931                     yylex();
 932 }
```
_____*unchanged_portion_omitted_*