Compiler Project01_Scanner Report

2015004693 양상헌

Goal:

Implementation of C-MINUS Scanner

Process:

- (1). By editing source codes "globals.h", "main.c", "util.h", "util.c", "scan.h", "scan.c"
- (2). By using flex by editing "tiny.I" to "cminus.I"
- (3). Creating both "cminus_cimpl", "cminus_flex" execution files by editing "Makefile"

(1).

일단 과제 명세 pdf에 나온 hint 대로 "globals.h" 에서 "MAXRESERVED"를 8에서 12로 바꿔주고, "TokenType"에 새로운 enum 변수들을 추가해 주었고, "main.c"에서 NO_PARSE, EchoSource, TraceScan을 TRUE로 바꿔주었고, "scan.c"에서 'StateType'과 'reservedWords[]'를 편집해주었다. 또한 "util.c" 에서 printToken()함수도 새로 추가된 tokentype들을 다 추가하여 편집해주었다.

가장 시간을 많이 할애하여 편집한 부분은 "scan.c"에서 'getToken()'함수의 내부 이다. 새로 추가된 StateType 에서의 경우의 수를 다 생각하여 DFA를 편집해주었다.

'=', '<', '>', '/', '!' 를 읽을 경우 각각 ASSIGN or EQ, LT or LE, GT or GE, OVER or INCOMMENT, NE or ERROR 가 될 수 있기에 이러한 경우를 잘 고려하여 DFA를 설계하였다. 또한 INCOMMENT상태에서 '*'를 읽게되면 INCOMMENT_ 상태가 되고 이때 '/'를 읽게되면 START, 다른 걸 일게 되면 다시 INCOMMENT 상태로 가게 된다.

편집된< getToken() in scan.c>, 빨간색 글자가 변경된 부분:

```
(1)./* function getToken returns the
                                                               (3)
                                                                          break;
 * next token in source file*/
TokenType getToken(void)
                                                                       case INEQ:
                                                                           //have read '='
  /* index for storing into tokenString */
  int tokenStringIndex = 0:
                                                                           state = DONE;
  /* holds current token to be returned */
  TokenType currentToken;
                                                                                      currentToken = E0:
  /* current state - always begins at START */
StateType state = START;
                                                                           else{
   /* flag to indicate save to tokenString */
                                                                                      ungetNextChar();
  int save;
                                                                                      save = FALSE:
  while (state != DONE)
                                                                                      currentToken = ASSIGN;
  { int c = getNextChar();
save = TRUE;
                                                                           break;
    switch (state)
                                                                        case INLT:
                                                                           //have read '<'
       case START:
       if (isdigit(c))
  state = INNUM;
                                                                           state = DONE:
                                                                                      currentToken = LE;
        else if (isalpha(c))
         state = INID;
                                                                           else{
                                                                                      ungetNextChar();
       else if ((c == ' ') || (c == '\t') || (c ==
                                                                                      save = FALSE:
'\n'))
                                                                                      currentToken = LT;
          save = FALSE;
                                                                           break;
         else if( c == '=' ){
                       save = FALSE;
                                                                        case INGT:
                                                                           //have read '>'
                       state = INEQ;
                                                                           state = DONE;
          else if( c =='<'){
                       save = FALSE;
                                                                                      currentToken = GE;
                       state = INLT
```

```
(2)
                                                                                (4)
               else if(c == '>'){
                                                                                                             ungetNextChar();
save = FALSE;
currentToken = GT;
                             save = FALSE;
                             state = INGT;
               break:
                             state = INNE;
                                                                                            case INNE:
                                                                                               //have read '!'
               else if(c == '/'){
                                                                                               state = DONE;
                             save = FALSE;
                                                                                               if(c == '='){
                                                                                                             currentToken = NE;
                             state = INOVER:
                                                                                               else{
            else
               state = DONE;
                                                                                                             ungetNextChar();
save = FALSE;
currentToken = ERROR;
              switch (c)
               { case EOF:
                   save = FALSE;
currentToken = ENDFILE;
                                                                                               break:
                                                                                            case INOVER:
                                                                                               // have read '/'
if(c == '*'){
                   case '+':
                   currentToken = PLUS;
                                                                                                             state = INCOMMENT;
save = FALSE;
                   currentToken = MINUS;
                                                                                               else {
                                                                                                             state = DONE:
                   break:
                                                                                                             ungetNextChar();
                   case '*':
currentToken = TIMES;
                                                                                                             save = FALSE;
                                                                                                             currentToken = OVER;
                   break:
                                                                                               break;
                   case '(':
currentToken = LPAREN;
                                                                                        case INCOMMENT:
                                                                                          save = FALSE;
if (c == EOF)
{ state = DONE;
                   break;
                    case ')':
                   currentToken = RPAREN;
                                                                                            currentToken = ENDFILE;
                                                                                          case '{':
                     currentToken = LCURLY;
                      break;
                                                                                          break;
                                                                                            case INCOMMENT_:
    // have read '/*---*'
    if(c == '/'){
                     case '}':
                      currentToken = RCURLY;
                      break;
                                                                                                             state = START;
                   case '[':
    currentToken = LBRACE;
                                                                                                             save = FALSE;
                                                                                               else if(c == '*'){
     state = INCOMMENT_;
     save = FALSE;
                      break:
                     case ']':
  currentToken = RBRACE;
                       break;
                                                                                               else if(c == EOF){
                                                                                                             state = DONE;
save = FALSE;
currentToken = ENDFILE;
                   case ';':
currentToken = SEMI;
                   break;
                                                                                               else{
                                                                                                             state = INCOMMENT;
save = FALSE;
                   case ',':
                     currentToken = COMMA;
                     break:
                                                                                               break;
                 currentToken = ERROR;
                                                                                            case DONE:
                 break;
                                                                                        default: /* should never happen */
fprintf(listing,"Scanner Bug:
                                                                                state= %d\n",state);
state = DONE;
          break;
       case INNUM:
                                                                                          currentToken = ERROR;
          if (!isdigit(c))
{ /* backup in the input */
  ungetNextChar();
                                                                                          break:
                                                                                      if ((save) && (tokenStringIndex <= MAXTOKENLEN))
  tokenString[tokenStringIndex++] = (char) c;</pre>
            save = FALSE;
state = DONE;
currentToken = NUM;
                                                                                          if (state == DONE)
                                                                                      { tokenString[tokenStringIndex] = '\0';
                                                                                        if (currentToken == ID)
  currentToken = reservedLookup(tokenString);
          break;
                                                                                     }
       case INID:
                                                                                   if (TraceScan) {
  fprintf(listing,"\t%d: ",lineno);
  printToken(currentToken,tokenString);
          if (!isalpha(c))
{ /* backup in the input */
  ungetNextChar();
            save = FALSE;
state = DONE;
currentToken = ID;
                                                                               return currentToken;
} /* end getToken */
```

(2).

"tiny,!"을 "cminus.!"로 편집하여 "flex"를 사용하여 "scan.c"의 역할을 하는 "lex,yy,c"를 생성한다. 이 부분에서 중요한 것은 아래 코드에서 빨간색으로 표시된 부분, NFA의 state를 나타내어 주는 부분을 잘 편집하는 것인 데. 'COMMENT'("/*...*/") state를 제외한 나머지 부분은 단순히 TokenType을 return 하는 방식으로 간단히 구

COMMENT(/*...*/)를 구현하는 부분에서 Nested comments(/* /* /* */ */)를 잘 인식할 수 있게 하는 부 분을 구현할 때 시간이 매우 많이 걸렸다.(이 부분은 밑에 cminus.l 에서 파란색 글씨로 되어있다)

"/*"를 읽고 INCOMMENT state로 들어왔을때, *를 읽으면 INCOMMENT_ state 가 되고 여기서 '/'을 읽으 면 START state가 되고, '/'이 아닌 다른것('\mathbb{'}n' '\mathbb{'}t' ' ' EOF 포함)을 읽게 되면 다시 INCOMMENT state 가 되게끔 만들어주었다. 이렇게 해줌으로써 '/*' 다음에 계속해서 '/*' 이 나오더라도 이를 단순 주석으로 인 식하고, 가장 먼저 나오는 '*/'로 INCOMMENT state를 벗어나게 된다. Ex("/* a /* b /* c */ d */ e */" 에서 "/* a /* b /* c */ d */ e */" 파란색 부분만 주석으로 인식이 되고 빨간색 부분은 일반 token Stream으로 인 식이 된다.)

편집된<cminus.l>, 빨간색 파란색 글자가 변경된 부분:

```
(1)
                                                                    (2)
                                                                                                         {return RBRACE:}
/* File: tiny.l
                                                                                                         {return LCURLY:}
                                                      */
/* Lex specification for TINY
                                                        */
                                                                                                         {return RCURLY;}
/* Compiler Construction: Principles and Practice
                                                                                     {return SEMI;}
/* Kenneth C. Louden
                                                                                                         {return COMMA;}
{number}
                                                                                      {return NUM:}
                                                                    {identifier}
                                                                                       {return ID;}
                                                                    {newline}
                                                                                      {lineno++;}
#include "globals.h"
#include "util.h"
#include "scan.h"
                                                                    {whitespace}
"/*" {
                                                                                      {/* skip whitespace */}
/* lexeme of identifier or reserved word */
                                                                                while(1){
char tokenString[MAXTOKENLEN+1];
                                                                                            while(c !='*'){
                                                                                            c = input();
if(c == EOF)
                                                                                             break;
if(c == '\n')
             [0-9]
diait
number
             {digit}+
             [a-zA-Z]
                                                                                               lineno++;
letter
identifier {letter}+
newline \n
whitespace [\t]+
                                                                                            //shoud be c == "*"
                                                                                            c = input();
if(c == EOF)
                                                                                            break;
if(c == '\n')
"if"
                                                                                            lineno++;
if(c == '/')
                 {return IF:}
"then"
                 {return THEN;}
                 {return ELSE;}
                                                                                               break;
"while"
                        {return WHILE;}
"return"
                        {return RETURN;}
"int"
                        {return INT;}
                                                                                    {return ERROR:}
"void"
                        {return VOID;}
"end"
                 {return END;}
{return REPEAT;}
"repeat"
                                                                    TokenType getToken(void)
"until"
                 {return UNTIL;}
"read"
                                                                    { static int firstTime = TRUE;
                 {return READ:}
"write"
                 {return WRITE;}
                                                                      TokenType currentToken;
                                                                      if (firstTime)
{ firstTime = FALSE;
0 \stackrel{\circ}{=} 0
                        {return ASSIGN;}
0==0
                 {return E0:}
^{n}\,!\!=^{n}
                        {return NE;}
                                                                        lineno++:
"<="
                        {return LE;}
                                                                        yyin = source;
">="
                        {return GE;}
                                                                        yyout = listing;
">"
                                    {return GT:}
                                                                      currentToken = yylex();
strncpy(tokenString,yytext,MAXTOKENLEN);
"<"
                {return LT:}
                {return PLUS;}
{return MINUS;}
n_{\pm}n
0 \triangle 0
                                                                      if (TraceScan) {
"*"
"/"
"("
                                                                        fprintf(listing,"\t%d: ",lineno);
                 {return TIMES;}
                 {return OVER;}
                                                                        printToken(currentToken, tokenString);
                {return LPAREN:}
")"
                {return RPAREN;}
                                                                      return currentToken;
                                    {return LBRACE;}
```

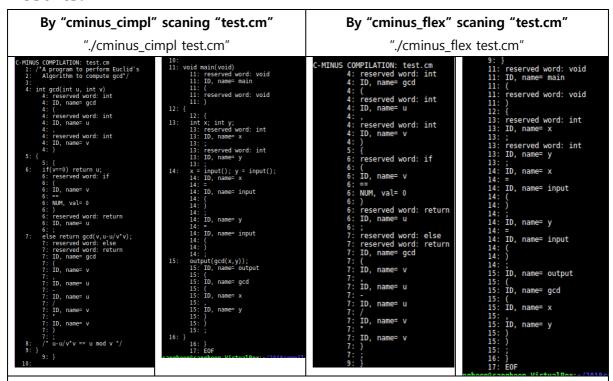
(3).

Makefile 에서 "make all" command를 이용하여 "cminus_cimpl", "cminus_flex" 실행파일들을 동시에 생성할수 있게끔 하기 위해 Makefile을 다음과 같이 편집해주었다.

편집된<Makefile>, 빨간색 글자가 변경된 부분:

```
CC = gcc
CFLAGS =
OBJS = main.o util.o scan.o parse.o symtab.o analyze.o code.o cgen.o
OBJS FLEX = main.o util.o lex.yy.o
cminus_cimpl: $(OBJS)
    $(CC) $(CFLAGS) $(OBJS) -o cminus_cimpl
main.o: main.c globals.h util.h scan.h parse.h analyze.h cgen.h
    $(CC) $(CFLAGS) -c main.c
util.o: util.c util.h globals.h $(CC) $(CFLAGS) -c util.c
lex.yy.o: cminus.l scan.h util.h globals.h
     flex cminus.l
    $(CC) $(CFLAGS) -c lex.yy.c -lfl
symtab.o: symtab.c symtab.h
    $(CC) $(CFLAGS) -c symtab.c
code.o: code.c code.h globals.h
    $(CC) $(CFLAGS) -c code.c
cgen.o: cgen.c globals.h symtab.h code.h cgen.h
    $(CC) $(CFLAGS) -c cgen.c
clean:
    -rm cminus_cimpl
                      #for scan.c
                      #for both
    -rm tm
    -rm cminus_flex #
-rm $(OBJS) #for both
                      #for flex
    $(CC) $(CFLAGS) tm.c -o tm
all: cminus_cimpl tm cminus_flex
```

Results:



과제 명세 pdf 내에 있는 예시 코드 test.cm을 scan 한 결과 위와 같은 출력 화면을 볼 수 있었다. 위 출력된 화면을 과제 명세 pdf 내에 있는 scan 결과 슬라이드와 비교해본 결과 모두 일치하였다.

- 또한 test.cm 의 주석 부분을 중첩으로 사용하여 실행해본 결과, 구현목표대로 "/* a /* b /* c */ d */ e */" 에서 "/* a /* b /* c */ d */ e */" 파란색 부분만 주석으로 인식이 되고 빨간색 부분은 일반 token Stream으로 인식이 되는 결과를 확인할 수 있었다.