

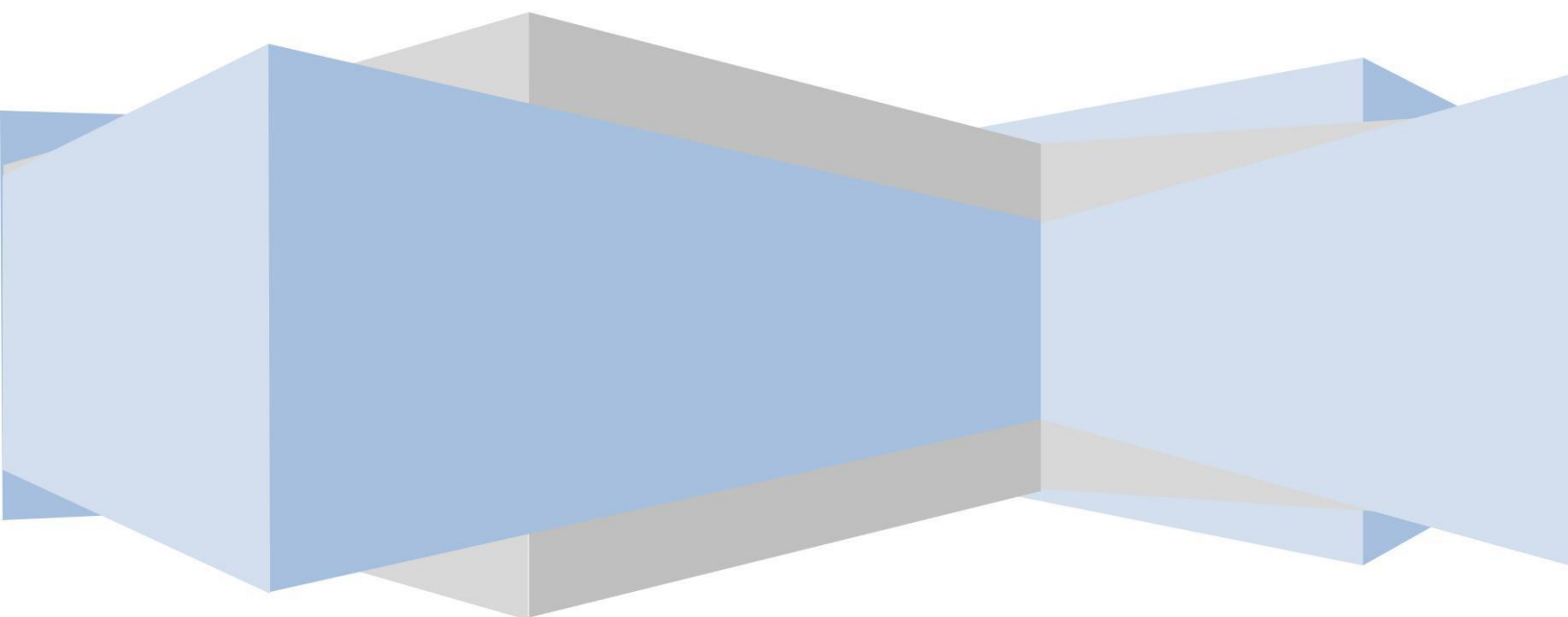
International Institute of Information Technology, Hyderabad

Compilers Project

Compiler for Decaf

Mohit Sharma

201505508



Objective:

Our objective for phase 1 is to implement the syntax analyzer and parser for the Decaf programming Language. Analyzer should be able to parse all valid programs and give an error for invalid program. Check [Decaf manual](#) to get an idea about the language.

Tools used:

1. **Flex:** Flex (fast lexical analyzer generator) is a free and open-source software alternative to lex. It is a computer program that generates lexical analyzers (also known as "scanners" or "lexers").

In simple words we can consider it as a tool which provides an easy way to specify the regular expressions corresponding to lexemes expected in source file. It also provides us a way to specify the action to be taken when we encounter a particular lexeme. In action part we can do operations like printing the lexeme (text matched) or returning corresponding token to the parser program. Check the following example. Here we have defined the regex for "if" and returning corresponding token in action part.

```
%{
    /* Global declarations and control information*/
}%

IF      "if"
%%
    {IF} { foutlex<<"IF"<<"\n";yylval.sval = strdup(yytext); return If; }
%%
```

The token we are returning here are specified in the parser file (bison file) explained below.

We pass this .l file to flex to generate corresponding C code for lexical analyzer which can scan any text for the lexemes specified in .l file. The output file from flex is named "lex.yy.c".

2. **Bison:** It is a parser generator that is part of the GNU Project. Bison reads a specification of a context-free language, warns about any parsing ambiguities, and generates a parser (either in C, C++, or Java) which reads sequences of tokens and decides whether the sequence conforms to the syntax specified by the grammar. Bison by default generates LALR parsers but can also create GLR parsers.

We specify the grammar and tokens that are needed for that grammar in a .y file. We provide this .y file to bison as input and it generates two files ".tab.h" (to be included in flex code, it contains all the terminal token declarations that flex will return when it encounters the lexeme corresponding to that token) and ".tab.c" (contains the parser code). To understand how to specify the grammar specifications, consider following dummy grammar which contains just one string "x y".

```
%{
    /* Global declarations and control information */
}%

// Union for data types of tokens that can be returned by flex
%union {
    int ival;
    char *sval;
}

// Different tokens that can be returned by flex
%token <sval> X Y

%%

    // Grammar that'll be parsed by bison
```

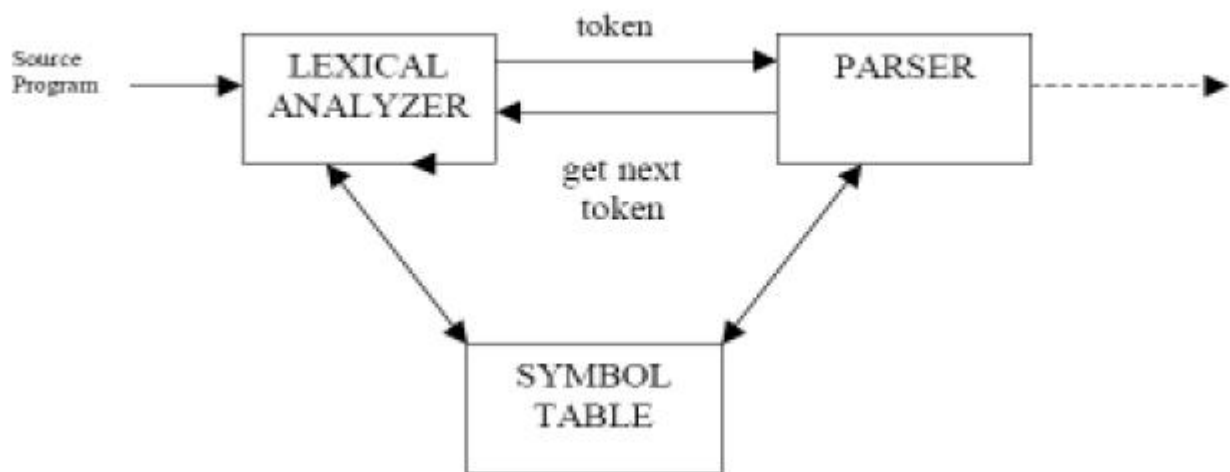
```

Program
    :    X Y    {fout <<"PROGRAM ENCOUNTERED\n"; };
%%

```

In this grammar we are expecting just two lexemes or terminals: “x” and “y”. Both are specified as tokens. The same are used in the only production rule. The parser generated will work in bottom up manner i.e. it'll first try to match X then Y, if it is able to complete the rule “X Y”, reduction will happen and matched rule will be replaced by corresponding LHS i.e. “Program” in this case. The moment it is able to reduce to start symbol, parser stops successfully. Whenever it reduces, the corresponding statements specified in the action part will be executed. If any unexpected symbol is encountered it calls the error handler “yyerror()”.

Check complete codes for flex and bison at the end.



Basic diagram to understand the interaction between lexical analyzer and parser.

Issues faced:

1. **Shift/reduce warnings:** These warnings can occur whenever grammar contains productions which can be reduced to a non-terminal and also one more symbol can be shifted instead of reducing it. If we don't specify the precedence and associativity of operators then there will be many shift/reduce warnings. The token specified earlier has lower precedence than that of specified later. For specifying associativity, we can use %left, %right and %noassoc. E.g. following three tokens are left associative.

```
%left Multiplication Division Remainder
```

2. **Same token to be used in two different contexts (binary/unary minus):** Both the minus are represented by '-' but we can't return two different tokens from flex corresponding to same string '-'. We need a way to use same token as per context. It can be done by specifying a dummy token for Unary minus just to specify that it has higher precedence than binary minus. Same is specified in the production using %prec. Check following example for better understanding:

```

...
%left '+' '-'
%left '*'
%left UMINUS

```

Now the precedence of UMINUS can be used in specific rules:

```
exp:
    ...
    | exp '-' exp
    ...
    | '-' exp %prec UMINUS
```

References:

1. [https://en.wikipedia.org/wiki/Flex_\(lexical_analyser_generator\)](https://en.wikipedia.org/wiki/Flex_(lexical_analyser_generator))
2. https://en.wikipedia.org/wiki/GNU_bison
3. http://www.gnu.org/software/bison/manual/html_node/Contextual-Precedence.html
4. <http://ashimg.tripod.com/Parser.html>
5. http://aquamentus.com/flex_bison.html

Codes

```
%{
/* Lex code to implement a parser for a decaf grammar */

#define YY_DECL extern "C" int yylex()
#include<fstream>
#include <iostream>
using namespace std;

#include "project.tab.h" // to get the token types that we return
fstream foutlex("flex_output.txt",ios::out);
%}

COMMENT          "//".*
CLASS             "class"
PROGRAM_STRING   "Program"
IDENTIFIER        [_a-zA-Z][_a-zA-Z0-9]*
OPEN_BRACKET      "{"
CLOSE_BRACKET     "}"
OPEN_SQUARE       "["
CLOSE_SQUARE      "]"
SEMICOLON         ";"
OPEN_PAREN        "("
CLOSE_PAREN       ")"
COMMA             ","
VOID              "void"
BREAK             "break"
CALLOUT           "callout"
CONTINUE          "continue"
ELSE              "else"
RETURN            "return"
EQUAL             "="
PLUS_EQUAL        "+="
MINUS_EQUAL       "-="
FOR               "for"
IF                "if"
INT               "int"
BOOLEAN           "boolean"
DECIMAL_LITERAL   [-]?[0-9]+
HEX_LITERAL       0[xX][0-9a-fA-F]+
CHAR_LITERAL      "\\'"([^\\"|'\\'|"\\\\"'|"\\\\\\'|"\\\\\\t"|"\\\\\\n")"'
STRING_LITERAL    "\".*\""
FALSE             "false"
TRUE              "true"
CONDITIONAL_OR    "||"
CONDITIONAL_AND   "&&"
EQUAL_TO          "=="
NOT_EQUAL_TO      "!="
LESS_THAN         "<"
LESS_THAN_EQUAL   "<="
GREATER_THAN_EQUAL ">="
GREATER_THAN      ">"
ADDITION          "+"
MINUS             "-"
MULTIPLICATION    "*"
DIVISION          "/"
REMAINDER         "%"
```

```

NEGATION                "!"

%%
[ \t\n]                 ;
{COMMENT}               ;

{CLASS}                  { foutlex<<"CLASS"<<"\n";
    yylval.sval = strdup(yytext); return Class; }
{PROGRAM_STRING}         { foutlex<<"PROGRAM"<<"\n";
    yylval.sval = strdup(yytext); return ProgramString; }
{VOID}                   { foutlex<<"VOID"<<"\n";
    yylval.sval = strdup(yytext); return Void; }
{BREAK}                  { foutlex<<"BREAK"<<"\n";
    yylval.sval = strdup(yytext); return Break; }
{CALLOUT}                { foutlex<<"CALLOUT"<<"\n";
    yylval.sval = strdup(yytext); return Callout; }
{CONTINUE}               { foutlex<<"CONTINUE"<<"\n";
    yylval.sval = strdup(yytext); return Continue; }
{ELSE}                   { foutlex<<"ELSE"<<"\n";
    yylval.sval = strdup(yytext); return Else; }
{RETURN}                 { foutlex<<"RETURN"<<"\n";
    yylval.sval = strdup(yytext); return Return; }
{FOR}                    { foutlex<<"FOR"<<"\n";
    yylval.sval = strdup(yytext); return For; }
{IF}                     { foutlex<<"IF"<<"\n";
    yylval.sval = strdup(yytext); return If; }
{INT}                    { foutlex<<"INT_DECLARATION"<<"\n";
    yylval.sval = strdup(yytext); return Int; }
{BOOLEAN}                { foutlex<<"BOOLEAN_DECLARATION"<<"\n";
    yylval.sval = strdup(yytext); return Boolean; }
{FALSE}                  { foutlex<<"BOOLEAN:false"<<"\n";
    yylval.sval = strdup(yytext); return False; }
{TRUE}                   { foutlex<<"BOOLEAN:true"<<"\n";
    yylval.sval = strdup(yytext); return True; }
{MINUS}                  { foutlex<<"-"<<"\n";
    yylval.sval = strdup(yytext); return Minus; }
{IDENTIFIER}             { foutlex<<"ID:"<<yytext<<"\n";
    yylval.sval = strdup(yytext); return Identifier; }
{OPEN_BRACKET}           { yylval.sval = strdup(yytext); return
OpenBracket; }
{CLOSE_BRACKET}          { yylval.sval = strdup(yytext); return
CloseBracket; }
{OPEN_SQUARE}            { yylval.sval = strdup(yytext); return
OpenSquare; }
{CLOSE_SQUARE}           { yylval.sval = strdup(yytext); return
CloseSquare; }
{SEMICOLON}              { yylval.sval = strdup(yytext); return
SemiColon; }
{OPEN_PAREN}             { yylval.sval = strdup(yytext); return
OpenParen; }
{CLOSE_PAREN}            { yylval.sval = strdup(yytext); return
CloseParen; }
{COMMA}                  { yylval.sval = strdup(yytext); return
Comma; }
{EQUAL}                  { foutlex<<"="<<"\n";
    yylval.sval = strdup(yytext); return Equal; }
{PLUS_EQUAL}             { foutlex<<"+="<<"\n";
    yylval.sval = strdup(yytext); return PlusEqual; }

```

```

{MINUS_EQUAL}          { foutlex<<"=="<<"\n";
                        yylval.sval = strdup(yytext); return MinusEqual; }
{DECIMAL_LITERAL}      { foutlex<<"INT:"<<yytext<<"\n";
                        yylval.sval = strdup(yytext); return Decimal_literal; }
{HEX_LITERAL}          { foutlex<<"INT:"<<yytext<<"\n";
                        yylval.sval = strdup(yytext); return Hex_literal; }
{CHAR_LITERAL}         { foutlex<<"CHARACTER:"<<yytext<<"\n";
                        yylval.sval = strdup(yytext); return Char_literal; }
{STRING_LITERAL}       { foutlex<<"STRING:"<<yytext<<"\n";
                        yylval.sval = strdup(yytext); return String_literal; }
{CONDITIONAL_OR}       { yylval.sval = strdup(yytext); return
ConditionalOr; }
{CONDITIONAL_AND}      { yylval.sval = strdup(yytext); return
ConditionalAnd; }
{EQUAL_TO}             { yylval.sval = strdup(yytext); return
EqualTo; }
{NOT_EQUAL_TO}         { yylval.sval = strdup(yytext); return
NotEqualTo; }
{LESS_THAN}           { yylval.sval = strdup(yytext); return
LessThan; }
{LESS_THAN_EQUAL}     { yylval.sval = strdup(yytext); return
LessThanEqual; }
{GREATER_THAN_EQUAL}  { yylval.sval = strdup(yytext); return
GreaterThanEqual; }
{GREATER_THAN}        { yylval.sval = strdup(yytext); return
GreaterThan; }
{ADDITION}            { yylval.sval = strdup(yytext); return
Addition; }
{MULTIPLICATION}      { yylval.sval = strdup(yytext); return
Multiplication; }
{DIVISION}            { yylval.sval = strdup(yytext); return
Division; }
{REMAINDER}           { yylval.sval = strdup(yytext); return
Remainder; }
{NEGATION}            { yylval.sval = strdup(yytext); return
Negation; }
%%

```

```

%{
// Bison code to parse decaf
#include <cstdio>
#include<string.h>
#include<fstream>
#include <iostream>
using namespace std;

// stuff from flex that bison needs to know about:
extern "C" int yylex();
extern "C" int yyparse();
extern "C" FILE *yyin;
void yyerror(const char *s);
extern "C" fstream foutlex;
fstream fout("bison_output.txt",ios::out);
%}
// Union for data types of tokens that can be returned by flex
%union {
    int ival;
    char *sval;
}

// Different tokens that can be returned by flex
%token <sval> Class ProgramString Identifier OpenBracket CloseBracket
OpenSquare CloseSquare SemiColon OpenParen CloseParen Comma Void
Break Callout Continue Else Return Equal PlusEqual MinusEqual For If
Int Boolean Decimal_literal Hex_literal Char_literal String_literal
True False
%left ConditionalOr
%left ConditionalAnd
%left EqualTo NotEqualTo
%left LessThan LessThanEqual GreaterThanEqual GreaterThan
%left Addition Minus
%left Multiplication Division Remainder
%left Negation
%left UMinus

%%
// Grammar that'll be parsed by bison
Program
    :      Class ProgramString OpenBracket field_decls method_decls
CloseBracket      {fout <<"PROGRAM ENCOUNTERED\n"; };

field_decls
    :      field_decls field_decl

    |;

field_decl
    :      Type Identifiers SemiColon;

Identifiers
    :      Identifier1
    |      Identifiers Comma Identifier1;

Identifier1

```



```

        :      Identifier
                                { fout<<"ID="<<$1<<"\n";
    }
    |      Identifier OpenSquare Decimal_literal CloseSquare
                                { fout<<"ID="<<$1<<"\n"<<"SIZE="<<$3<<"\n"; }
    |      Identifier OpenSquare Hex_literal CloseSquare
                                {
fout<<"ID="<<$1<<"\n"<<"SIZE="<<$3<<"\n"; };
method_decls
    :      method_decl method_decls

        |;

method_decl
    :      Type Identifier OpenParen Params CloseParen Block
                                { fout<<"METHOD="<<$2<<"\n"; }
    |      Type Identifier OpenParen CloseParen Block
                                { fout<<"METHOD="<<$2<<"\n"; }
    |      Void Identifier OpenParen Params CloseParen Block
                                { fout<<"METHOD="<<$2<<"\n"; }
    |      Void Identifier OpenParen CloseParen Block
                                { fout<<"METHOD="<<$2<<"\n"; };

Params
    :      Type Identifier

    |      Type Identifier Comma Params;

Block
    :      OpenBracket var_decls Statements CloseBracket;

var_decls
    :      var_decl var_decls

    |;

var_decl
    :      Type decls SemiColon;

decls
    :      Identifier
                                { fout<<"ID="<<$1<<"\n"; }
    |      Identifier Comma decls
                                { fout<<"ID="<<$1<<"\n"; };

Statements
    :      Statement Statements

    |;

Statement
    :      Location Assign_op Expr SemiColon
                                { fout<<"ASSIGNMENT OPERATION
ENCOUNTERED\n"; }
    |      Method_call SemiColon
    |      If OpenParen Expr CloseParen Block
                                { fout<<"IF ENCOUNTERED\n"; }

```

```

|      If OpenParen Expr CloseParen Block Else Block
|                                { fout<<"IF ENCOUNTERED\n"; }
|      For Identifier Equal Expr Comma Expr Block
|                                { fout<<"FOR ENCOUNTERED\n"; }
|      Return SemiColon
|                                { fout<<"RETURN ENCOUNTERED\n"; }
|      Return Expr SemiColon
|                                { fout<<"RETURN ENCOUNTERED\n"; }

|      Break SemiColon
|                                { fout<<"BREAK ENCOUNTERED\n";
}

|      Continue SemiColon
|                                { fout<<"CONTINUE
ENCOUNTERED\n"; }
|      Block;

Method_call
:      Method_name OpenParen PassParams CloseParen

|      Method_name OpenParen CloseParen
|      Callout      OpenParen String_literal CloseParen
|                                { fout<<"CALLOUT TO "<<$3<<"
ENCOUNTERED\n"; }
|      Callout      OpenParen String_literal Comma CalloutArgs
CloseParen
|                                { fout<<"CALLOUT TO "<<$3<<"
ENCOUNTERED\n"; };

PassParams
:      Expr
|      Expr Comma PassParams;

Method_name
:      Identifier
|                                { fout<<"METHOD CALL="<<$1; };

Location
:      Identifier
|                                { fout<<"LOCATION
ENCOUNTERED="<<$1<<"\n"; }
|      Identifier OpenSquare Expr CloseSquare;

CalloutArgs
:      Expr
|      String_literal
|      Expr Comma CalloutArgs
|      String_literal Comma CalloutArgs;

Expr
:      Location
|      Method_call
|      Literal
|      Expr ConditionalOr Expr
|                                { fout<<"CONDITIONAL OR
ENCOUNTERED\n"; }
|      Expr ConditionalAnd Expr
|                                { fout<<"CONDITIONAL AND
ENCOUNTERED\n"; }

```

```

        |      Expr EqualTo Expr
                                { fout<<"EQUAL TO
ENCOUNTERED\n"; }
        |      Expr NotEqualTo Expr
                                { fout<<"NOT EQUAL TO
ENCOUNTERED\n"; }
        |      Expr LessThan Expr
                                { fout<<"LESS THAN
ENCOUNTERED\n"; }
        |      Expr LessThanEqual Expr
                                { fout<<"LESS THAN EQUAL
ENCOUNTERED\n"; }
        |      Expr GreaterThanEqual Expr
                                { fout<<"GREATER THAN EQUAL
ENCOUNTERED\n"; }
        |      Expr GreaterThan Expr
                                { fout<<"GREATER THAN
ENCOUNTERED\n"; }
        |      Expr Addition Expr
                                { fout<<"ADDITION
ENCOUNTERED\n"; }
        |      Expr Minus Expr
                                { fout<<"SUBTRACTION
ENCOUNTERED\n"; }
        |      Expr Multiplication Expr
                                { fout<<"MULTIPLICATION
ENCOUNTERED\n"; }
        |      Expr Division Expr
                                { fout<<"DIVIDION
ENCOUNTERED\n"; }
        |      Expr Remainder Expr
                                { fout<<"MOD ENCOUNTERED\n"; }
        |      Minus Expr %prec UMinus
                                { fout<<"UNARY MINUS
ENCOUNTERED\n"; }
        |      Negation Expr
                                { fout<<"NEGATION
ENCOUNTERED\n"; }
        |      OpenParen Expr CloseParen;

Literal
    :      Int_literal
    |      Char_literal
                                { fout<<"CHAR
ENCOUNTERED=$1<<"\n"; }
    |      Bool_literal;

Bool_literal
    :      True
                                { fout<<"BOOLEAN
ENCOUNTERED=$1<<"\n"; }
    |      False
                                { fout<<"BOOLEAN
ENCOUNTERED=$1<<"\n"; };

Int_literal

```

```

        :      Decimal_literal
                                { fout<<"INT
ENCOUNTERED="<<$1<<"\n"; }
        |      Hex_literal
                                { fout<<"INT
ENCOUNTERED="<<$1<<"\n"; };

Type
    :      Int
                                { fout<<"INT DECLARATION
ENCOUNTERED\n"; }
        |      Boolean
                                { fout<<"BOOLEAN
DECLARATION ENCOUNTERED\n"; };

Assign_op
    :      Equal
                                { fout<<"ASSIGNMENT
ENCOUNTERED\n"; }
        |      PlusEqual
                                { fout<<"ADDITION ASSIGNMENT
ENCOUNTERED\n"; }
        |      MinusEqual
                                { fout<<"SUBTRACTION
ASSIGNMENT ENCOUNTERED\n"; };
%%

int main(int argc, char** argv)
{
    FILE *myfile = fopen(argv[1], "r");
    if (!myfile)
    {
        cout << "I can't open input file!" << endl;
        return -1;
    }
    yyin = myfile;
    do
    {
        yyparse();
    } while (!feof(yyin));
    cout<<"Success\n";
    fout.close();
    foutlex.close();
    return 0;
}

void yyerror(const char *s) {
    cout<<"Syntax error\n";
    fout.close();
    exit(-1);
}

```

```
# Use this file to run above codes
# Input: "test_input" containing decaf code
# Output: flex_output.txt, bison_output.txt
```

```
#!/bin/bash
```

```
bison -dv Project.y
flex Project.l
g++ Project.tab.c lex.yy.c -lfl -o proj
./proj test_input
```