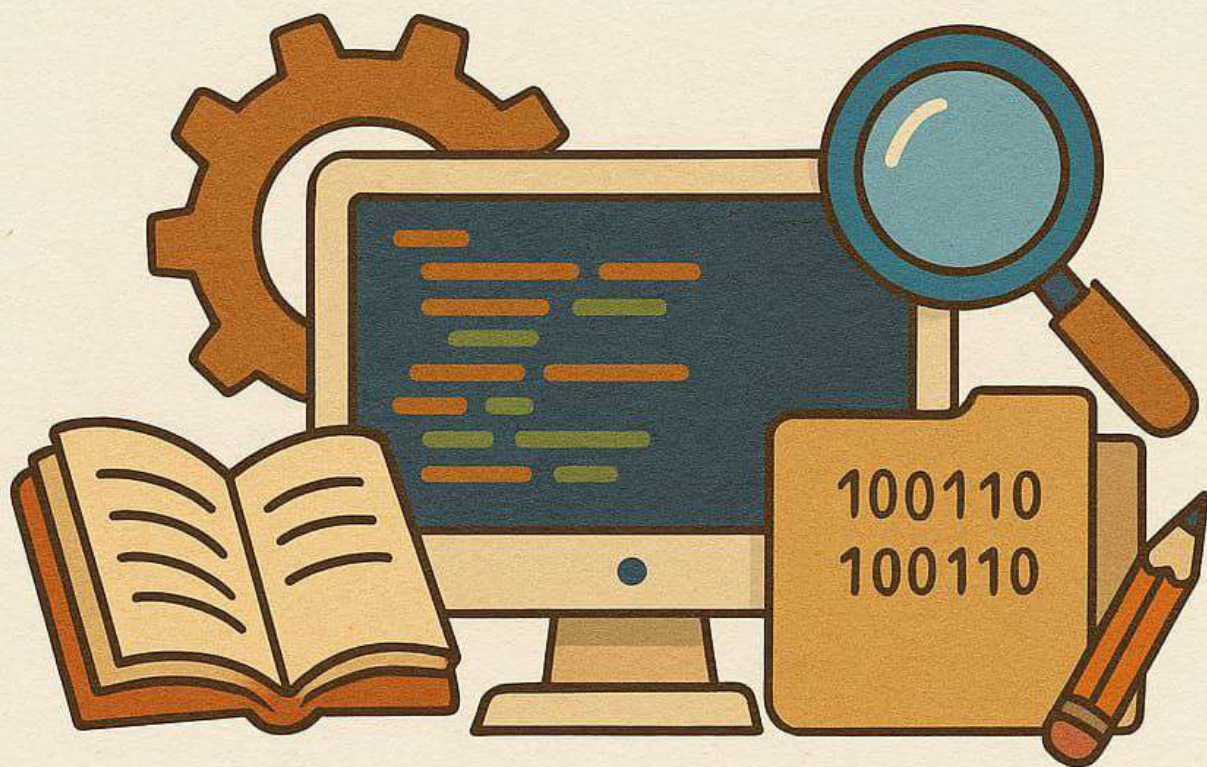


COMPILER DESIGN



Members:

Abhinandan Barua

Shahir Habib

Masud Alam

Sanket Chakraborty

Abstract

This report details the development of a compiler for a subset of the C programming language. The project encompasses the implementation of key compiler components: lexical analysis, syntax analysis using LR(1) parsing, symbol table management, and the generation of intermediate representations. The report outlines the methodologies employed, the structure of the implementation, challenges encountered, and potential areas for future enhancement.

Table of Contents

1. Introduction
2. Project Objectives
3. Implementation Details
 - 3.1 Lexical Analysis
 - 3.2 Syntax Analysis
 - 3.3 Symbol Table Management
 - 3.4 Intermediate Representations
4. Testing and Results
5. Challenges and Solutions
6. Future Work
7. Conclusion
8. References

Introduction

Compilers are essential tools that translate high-level programming languages into machine code, enabling program execution on hardware. This project focuses on constructing a compiler for a subset of the C language, implementing core components to understand the intricacies of compiler design and construction.

Project Objectives

The primary objectives of this project are:

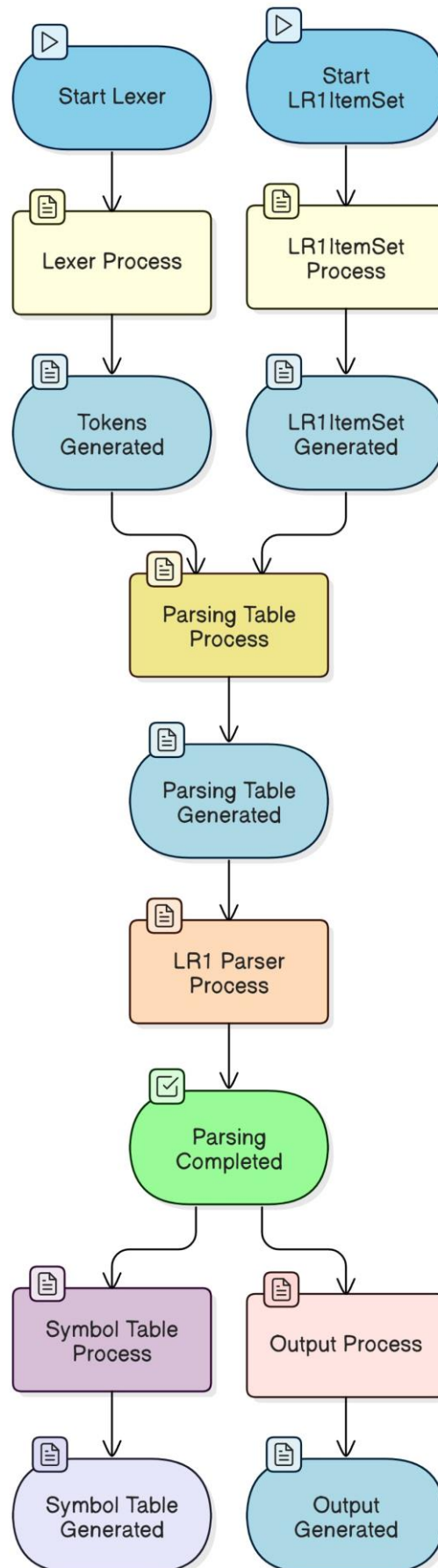
- To develop a functional compiler capable of processing a subset of the C language.
- To implement and integrate fundamental compiler components: lexical analysis, syntax analysis using LR(1) parsing, symbol table management, and intermediate representation generation.
- To gain practical experience with compiler construction techniques and tools.



Language Specification

- Data Types: int, float, void
- Input/Output: read, print
- Operators: +, -, *, /, %, ++
- Conditionals: if, else (no nesting)
- Relational Operators: <, >, ==
- Multiple functions allowed
- Declarations can appear anywhere with proper scope handling

Flow Diagram:



Implementation Details

The compiler development is structured into four main components: lexical analysis, syntax analysis, symbol table management, and Parsing.

File Name	Description
lexer.cpp	Performs lexical analysis and generates tokens
program.txt	Sample input C-like code
tokens.txt	Output of lexer containing tokens
symbolTable.cpp	Manages variables, types, scopes
symbol_table.txt	Stores final symbol table output
cfg.txt	Context-Free Grammar definition
lr1itemset.cpp	Constructs LR(1) item sets from grammar
lr1itemsets.txt	Output of LR(1) item sets
parsingTable.cpp	Generates LR(1) parsing table
parsingtable.txt	Stores action/goto parsing table
lr1parser.cpp	Uses parsing table to parse token stream
parsingresult.txt	Output showing successful parsing or error

Lexical Analysis

Objective: To scan the source code and convert it into a stream of tokens.

Implementation:

- Lexer Development (lexer.cpp): The lexer reads the source code from program.txt, identifies lexemes, and classifies them into tokens such as keywords, identifiers, literals, operators, and punctuation. The resulting tokens are output to tokens.txt.

Challenges:

- Accurately distinguishing between similar lexemes (e.g., distinguishing between the keyword if and an identifier ifelse).
- Efficiently handling white spaces and comments.

Solutions:

- Implemented a longest match rule to resolve ambiguities between similar lexemes.
- Designed the lexer to ignore white spaces and comments, ensuring they do not interfere with tokenization.

Context-Free Grammar (CFG):

Context Free Grammar Design

TOKENS ID NUM DEC INT VOID FLOAT EQ LPAREN RPAREN LBRACE RBRACE SEMI COMMA IF ELSE READ PRINT EQEQ LT GT PLUS MINUS MULT DIV MOD INC

```
S: function_list S
  | declaration_list S
  | function_list
  | declaration_list
  ;
```

```
declaration_list: declaration
  | declaration_list declaration
  ;
```

```
declaration: type_specifier ID SEMI
  ;
```

```
type_specifier: INT
  | VOID
  | FLOAT
  ;
```

```
function_list: function
  | function_list function
  ;
```

```
function: type_specifier ID LPAREN params
  RPAREN compound_stmt
  ;
```

```
params:
  | type_specifier ID COMMA params
  | type_specifier ID
  ;
```

```
compound_stmt: LBRACE statement_list
  RBRACE
  ;
```

```
statement_list: statement
  | statement_list statement
  ;
```

```
statement: declaration
  | assignments
  | return_stmt
  | condition
  | read_stmt
  | print_stmt
  ;
```

```
return_stmt: ID expression SEMI
  ;
```

```
assignments: ID EQ expression SEMI
  ;
```

```
expression: ID
  | NUM
  | DEC
  | expression PLUS expression
  | expression MINUS expression
  | expression MULT expression
  | expression DIV expression
  | expression MOD expression
  | expression INC
  | LPAREN expression RPAREN
  ;
```

```
condition: IF LPAREN rel_expression RPAREN
  compound_stmt ELSE compound_stmt
  ;
```

```
rel_expression: expression EQEQ expression
  | expression LT expression
  | expression GT expression
  ;
```

```
read_stmt: READ ID SEMI
  ;
```

```
print_stmt: PRINT expression SEMI
  ;
```

- Start Symbol: S defines a program that can consist of multiple declarations and functions.
- Declaration and Function Blocks:
 - *declaration_list* represents multiple variable declarations.
 - *function_list* allows defining one or more functions.
 - function has structure: return type + function name + parameters + compound statement.
- Parameter Parsing:

- *params* allows functions to have typed parameters or be empty.
- Compound Statement:
 - Block enclosed in {} that holds a list of valid statements.
- **Statements include:**
 - *declaration, assignments, return_stmt, condition, read_stmt, print_stmt*
- **Expressions:**
 - Support arithmetic: *+, -, *, /, %, and ++*
 - Can include variables, numbers, or nested expressions
- **Conditionals:**
 - Format: *if (rel_expression) { ... } else { ... }*
 - No nested if-else support
- **Relational Expressions:**
 - Comparisons using *==, <, and >*
- **I/O Statements:**
 - *read_stmt* and *print_stmt* use READ and PRINT tokens respectively

LR(1) Item Set Generation:

1. Reads the CFG.
2. Constructs LR(1) item sets using closure and goto functions.
3. Implements FIRST and FOLLOW computation.
4. Generates lr1itemsets.txt.

```
46      INC
47      _____P R O D U C T I O N S _____
48      0 S' -> S
49      1 S -> function_list S
50      2 S -> declaration_list S
51      3 S -> function_list
52      4 S -> declaration_list
53      5 declaration_list -> declaration
54      6 declaration_list -> declaration_list declaration
55      7 declaration -> type_specifier ID SEMI
56      8 type_specifier -> INT
57      9 type_specifier -> VOID
58      10 type_specifier -> FLOAT
59      11 function_list -> function
60      12 function_list -> function_list function
61      13 function -> type_specifier ID LPAREN params RPAREN compound_stmt
62      14 params ->
63      15 params -> type_specifier ID COMMA params
64      16 params -> type_specifier ID
65      17 compound_stmt -> LBRACE statement_list RBRACE
66      18 statement_list -> statement
67      19 statement_list -> statement_list statement
68      20 statement -> declaration
69      21 statement -> assignments
70      22 statement -> return_stmt
```

Follow:

```
138      S : $
139      S' :
140      assignments : FLOAT ID IF INT PRINT RBRACE READ VOID
141      compound_stmt : $ ELSE FLOAT ID IF INT PRINT RBRACE READ VOID
142      condition : FLOAT ID IF INT PRINT RBRACE READ VOID
143      declaration : $ FLOAT ID IF INT PRINT RBRACE READ VOID
144      declaration_list : $ FLOAT INT VOID
145      expression : DIV EQEQ GT INC LT MINUS MOD MULT PLUS RPAREN SEMI
146      function : $ FLOAT INT VOID
147      function_list : $ FLOAT INT VOID
148      params : RPAREN
149      print_stmt : FLOAT ID IF INT PRINT RBRACE READ VOID
150      read_stmt : FLOAT ID IF INT PRINT RBRACE READ VOID
151      rel_expression : RPAREN
152      return_stmt : FLOAT ID IF INT PRINT RBRACE READ VOID
153      statement : FLOAT ID IF INT PRINT RBRACE READ VOID
154      statement_list : FLOAT ID IF INT PRINT RBRACE READ VOID
155      type_specifier : ID
```

First:

```
117      S : FLOAT INT VOID
118      S' : FLOAT INT VOID
119      SEMI : SEMI
120      VOID : VOID
121      assignments : ID
122      compound_stmt : LBRACE
123      condition : IF
124      declaration : FLOAT INT VOID
125      declaration_list : FLOAT INT VOID
126      expression : DEC ID LPAREN NUM
127      function : FLOAT INT VOID
128      function_list : FLOAT INT VOID
129      params : FLOAT INT VOID epsilon
130      print_stmt : PRINT
131      read_stmt : READ
132      rel_expression : DEC ID LPAREN NUM
133      return_stmt : ID
134      statement : FLOAT ID IF INT PRINT READ VOID
135      statement_list : FLOAT ID IF INT PRINT READ VOID
136      type_specifier : FLOAT INT VOID
```

Parsing Table Generation:

1. Uses the LR(1) item sets to compute the parsing table.
2. Determines shift/reduce actions and goto transitions.
3. Outputs parsingtable.txt.

----- LR(1) Parsing Table -----									
-----+-----									
MINUS	MULT	DIV	MOD	INC	\$	S	declaration_list	dec	
						3	6		
					Accept				
					r5				
					r4	10	6		
					r11				
					r3	13	6		
<									>

LR(1) Parsing:

1. Reads tokens.txt and parsingtable.txt.
2. Simulates the parsing stack.
3. Logs reductions, matches, and errors.

```

1      Lookahead : FLOAT
2          | Action : s1
3      Shift : FLOAT
4      Lookahead : ID
5          | Action : r10
6      Reduce by type_specifier -> FLOAT
7      Lookahead : ID
8          | Action : s16
9      Shift : ID
10     Lookahead : LPAREN
11         | Action : s19
12     Shift : LPAREN
13     Lookahead : RPAREN
14         | Action : r14
15     Reduce by params ->
16     Lookahead : RPAREN
17         | Action : s23
18     Shift : RPAREN
19     Lookahead : LBRACE
20         | Action : s25
21     Shift : LBRACE
22     Lookahead : ID
23         | Action : s28
24     Shift : ID
25     Lookahead : DEC

```

```

158         | Action : r15
159     Reduce by statement_list -> statement_list statement
160     Lookahead : RBRACE
161         | Action : s51
162     Shift : RBRACE
163     Lookahead : $
164         | Action : r17
165     Reduce by compound_stmt -> LBRACE statement_list RBRACE
166     Lookahead : $
167         | Action : r13
168     Reduce by function -> type_specifier ID LPAREN params RPAREN compound_stmt
169     Lookahead : $
170         | Action : r11
171     Reduce by function_list -> function
172     Lookahead : $
173         | Action : r3
174     Reduce by S -> function_list
175     Lookahead : $
176         | Action : r2
177     Reduce by S -> declaration_list S
178     Lookahead : $
179         | Action : r1
180     Reduce by S -> function_list S
181     Lookahead : $
182         | Accepted 🍀 🍀 🍀 🍀

```


Symbol Table:

symbolTable.cpp

- Creates a symbol table during parsing.
- Stores variable names, types, and scope.

1	+-----+-----+-----+-----+-----+-----+-----+-----+							
2	Token Type	Name	Value	Line	Pos	Scope	Memory Addr	
3	+-----+-----+-----+-----+-----+-----+-----+-----+							
4	FLOAT	fun		1	0	0	0x1000	
5	+-----+-----+-----+-----+-----+-----+-----+-----+							
6	Referenced at:							
7								
8	→ Line 1							
9	+-----+-----+-----+-----+-----+-----+-----+-----+							
10	INT	x		3	4	1	0x1008	
11	+-----+-----+-----+-----+-----+-----+-----+-----+							
12	Referenced at:							
13								
14	→ Line 3							
15	+-----+-----+-----+-----+-----+-----+-----+-----+							
16	INT	x	55	7	5	0	0x100c	
17	+-----+-----+-----+-----+-----+-----+-----+-----+							
18	Referenced at:							
19								
20	→ Line 5							
21	→ Line 7							
22	+-----+-----+-----+-----+-----+-----+-----+-----+							
23	INT	main		6	0	0	0x1010	

Conclusion:

The project demonstrates the end-to-end design of a simplified compiler. It uses LR(1) parsing, proper tokenization, and symbol table management. The modular design enables future extensions like semantic analysis and code generation.

