

Bab 1

Implementasi Lexer Sederhana (Hand-Written)

1.1 Tujuan Pembelajaran

Setelah mempelajari bab ini, mahasiswa diharapkan mampu:

1. Memahami konsep dan struktur hand-written lexer
2. Merancang state machine untuk token recognition
3. Mengimplementasikan lexer sederhana dalam C/C++ untuk subset bahasa C
4. Menangani whitespace, komentar (single-line dan multi-line), dan escape sequences
5. Mengimplementasikan error handling untuk token tidak valid
6. Membuat unit test untuk berbagai kasus input

1.2 Pendahuluan

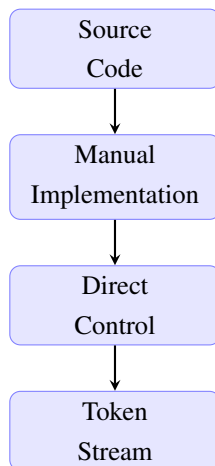
Setelah memahami konsep lexical analysis secara teori pada bab sebelumnya, pada bab ini kita akan mengimplementasikan lexer secara praktis menggunakan pendekatan **hand-written** (ditulis manual). Menurut sumber terbuka:

“Hand-written lexers are possible: directly code a state machine, or use manual scanning logic. Requires careful handling of edge cases (e.g. unclosed strings/comments).”?

Pendekatan hand-written memberikan kontrol penuh terhadap implementasi dan sangat berguna untuk pembelajaran karena mahasiswa dapat memahami setiap detail proses tokenization. Meskipun lebih kompleks dibanding menggunakan generator seperti Flex atau re2c, hand-written lexer memberikan fleksibilitas dan pemahaman yang lebih dalam.

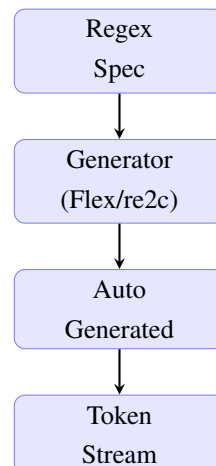
Gambar [1.1](#) menunjukkan perbandingan antara hand-written lexer dan lexer generator.

HAND-WRITTEN LEXER



Pros: Full control,
learning value
Cons: More code

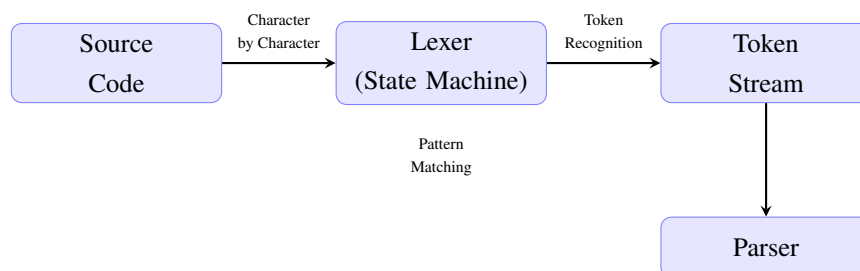
LEXER GENERATOR



Pros: Less code,
maintainable
Cons: Less control

Gambar 1.1: Perbandingan hand-written lexer vs lexer generator

Gambar 1.2 menunjukkan alur umum proses tokenization dalam hand-written lexer.



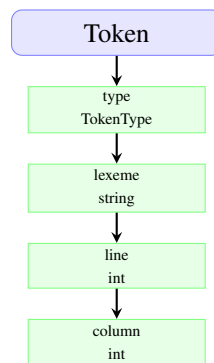
Gambar 1.2: Alur umum proses tokenization dalam hand-written lexer

1.3 Struktur Token

Sebelum mengimplementasikan lexer, kita perlu mendefinisikan struktur data untuk merepresentasikan token. Token minimal harus menyimpan:

1. **Token Type:** Jenis token (identifier, keyword, number, operator, dll.)
2. **Lexeme:** String aktual yang di-match dari source code
3. **Position Information:** Baris dan kolom untuk error reporting
4. **Value** (opsional): Nilai numerik untuk number literals

Gambar 1.3 menunjukkan struktur data token secara visual.

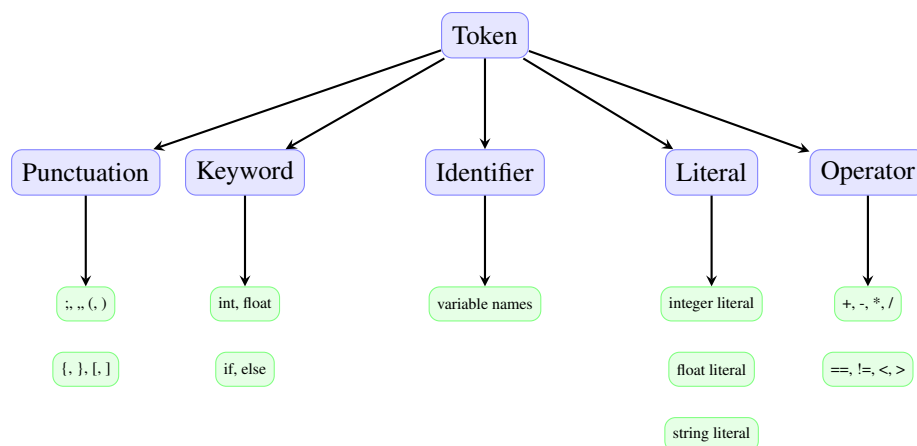


Gambar 1.3: Struktur data Token

1.3.1 Token Types

Token types dapat didefinisikan menggunakan enum. Berikut contoh untuk subset bahasa C:

Gambar 1.4 menunjukkan hierarki token types yang digunakan dalam lexer.



Gambar 1.4: Hierarki tipe token dalam lexer

Listing 1.1: Definisi Token Types

```

1 enum class TokenType {
2     // Identifiers and Keywords
3     IDENTIFIER,
4     KEYWORD_INT, KEYWORD_FLOAT, KEYWORD_IF, KEYWORD_ELSE,
5     KEYWORD_WHILE, KEYWORD_FOR, KEYWORD_RETURN,
6
7     // Literals
8     INTEGER_LITERAL,
9     FLOAT_LITERAL,
10    STRING_LITERAL,
11    CHAR_LITERAL,

```

```

12
13 // Operators
14 OP_PLUS, OP_MINUS, OP_MULTIPLY, OP_DIVIDE,
15 OP_ASSIGN, OP_EQUAL, OP_NOT_EQUAL,
16 OP_LESS, OP_LESS_EQUAL, OP_GREATER, OP_GREATER_EQUAL,
17 OP_AND, OP_OR, OP_NOT,
18
19 // Punctuation
20 SEMICOLON, COMMA, DOT,
21 LPAREN, RPAREN, // ( )
22 LBRACE, RBRACE, // { }
23 LBRACKET, RBRACKET, // [ ]
24
25 // Special
26 END_OF_FILE,
27 INVALID
28 };

```

1.3.2 Token Structure

Struktur token dalam C++ dapat didefinisikan sebagai berikut:

Listing 1.2: Struktur Token

```

1 struct Token {
2     TokenType type;
3     std::string lexeme;
4     int line;
5     int column;
6     union {
7         int intValue; // Untuk INTEGER_LITERAL
8         double floatValue; // Untuk FLOAT_LITERAL
9     };
10
11     Token(TokenType t, const std::string& lex, int l, int c)
12         : type(t), lexeme(lex), line(l), column(c) {}
13 };

```

1.4 Finite State Machine untuk Lexer

Lexical analysis secara fundamental adalah proses pattern matching yang dapat dimodelkan menggunakan **Finite State Machine (FSM)** atau **Finite Automata**. Menurut sumber dari Aoyama Gakuin University:

“Lexical analysis breaks input text into lexemes which correspond to tokens. Usually implemented using regular languages → regex → NFA → DFA → (minimized) DFA for efficiency.”?

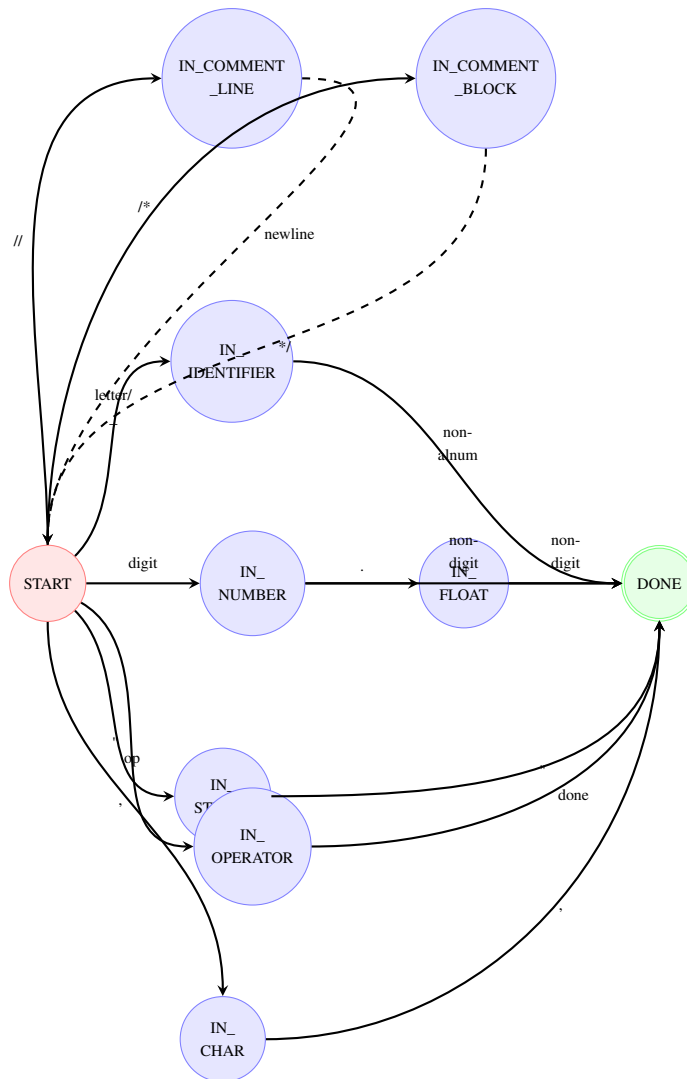
Dalam implementasi hand-written, kita tidak perlu membuat DFA secara eksplisit, tetapi kita menggunakan logika state machine dalam kode.

1.4.1 State Machine Design

State machine untuk lexer sederhana dapat memiliki state-state berikut:

- **START**: State awal, menunggu karakter pertama dari token
- **IN_IDENTIFIER**: Sedang membaca identifier atau keyword
- **IN_NUMBER**: Sedang membaca angka (integer atau float)
- **IN_FLOAT**: Setelah menemukan titik desimal
- **IN_STRING**: Sedang membaca string literal
- **IN_CHAR**: Sedang membaca character literal
- **IN_COMMENT_LINE**: Sedang membaca single-line comment
- **IN_COMMENT_BLOCK**: Sedang membaca multi-line comment
- **IN_OPERATOR**: Sedang membaca operator (mungkin multi-character)
- **DONE**: Token selesai dibaca

Gambar [1.5](#) menunjukkan state machine untuk lexer sederhana.



Gambar 1.5: State machine untuk hand-written lexer

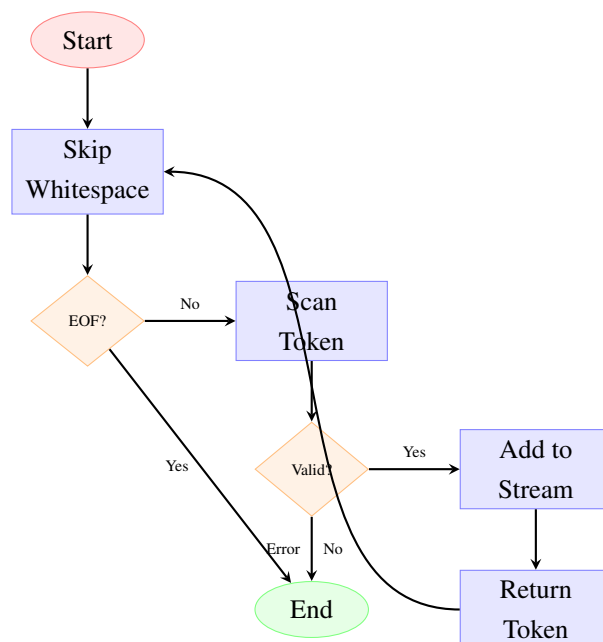
1.4.2 State Transitions

Transisi state terjadi berdasarkan karakter yang dibaca:

1. **START** → **IN_IDENTIFIER**: Jika karakter adalah huruf atau underscore
2. **START** → **IN_NUMBER**: Jika karakter adalah digit
3. **START** → **IN_STRING**: Jika karakter adalah double quote (")
4. **START** → **IN_CHAR**: Jika karakter adalah single quote (')
5. **START** → **IN_COMMENT_LINE**: Jika menemukan //
6. **START** → **IN_COMMENT_BLOCK**: Jika menemukan /*

7. **START** → **IN_OPERATOR**: Jika karakter adalah operator
8. **IN_NUMBER** → **IN_FLOAT**: Jika menemukan titik desimal
9. **IN_IDENTIFIER** → **DONE**: Jika karakter bukan alphanumeric atau underscore
10. **IN_NUMBER** → **DONE**: Jika karakter bukan digit atau titik
11. **IN_STRING** → **DONE**: Jika menemukan closing quote (dengan handling escape)

Gambar 1.6 menunjukkan flowchart proses tokenization.



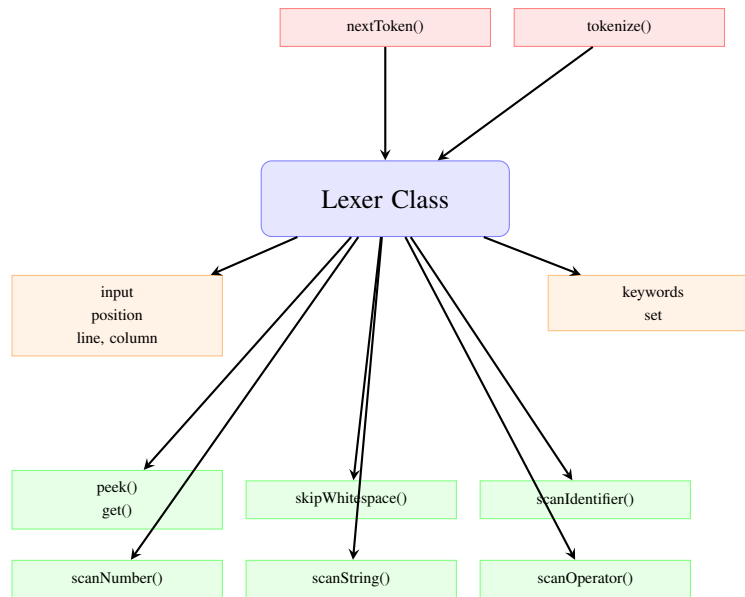
Gambar 1.6: Flowchart proses tokenization

1.5 Implementasi Lexer dalam C++

Berikut adalah implementasi lengkap lexer sederhana untuk subset bahasa C:

1.5.1 Kelas Lexer

Gambar 1.7 menunjukkan arsitektur kelas Lexer dan komponen-komponennya.



Gambar 1.7: Arsitektur kelas Lexer

Listing 1.3: Header File: lexer.h

```

1 #ifndef LEXER_H
2 #define LEXER_H
3
4 #include <string>
5 #include <unordered_set>
6 #include <vector>
7
8 enum class TokenType {
9     IDENTIFIER,
10    KEYWORD_INT, KEYWORD_FLOAT, KEYWORD_IF, KEYWORD_ELSE,
11    KEYWORD_WHILE, KEYWORD_FOR, KEYWORD_RETURN,
12    INTEGER_LITERAL, FLOAT_LITERAL,
13    STRING_LITERAL, CHAR_LITERAL,
14    OP_PLUS, OP_MINUS, OP_MULTIPLY, OP_DIVIDE,
15    OP_ASSIGN, OP_EQUAL, OP_NOT_EQUAL,
16    OP_LESS, OP_LESS_EQUAL, OP_GREATER, OP_GREATER_EQUAL,
17    OP_AND, OP_OR, OP_NOT,
18    SEMICOLON, COMMA, DOT,
19    LPAREN, RPAREN, LBRACE, RBRACE,
20    LBRACKET, RBRACKET,
21    END_OF_FILE, INVALID
22 };
23
24 struct Token {
25     TokenType type;
26     std::string lexeme;
27     int line;
28     int column;
29
30     Token(TokenType t, const std::string& lex, int l, int c)
31         : type(t), lexeme(lex), line(l), column(c) {}

```



```

32 };
33
34 class Lexer {
35 private:
36     std::string input;
37     size_t position;
38     int line;
39     int column;
40     std::unordered_set<std::string> keywords;
41
42     char peek() const;
43     char get();
44     void skipWhitespace();
45     void skipLineComment();
46     void skipBlockComment();
47     Token scanIdentifier();
48     Token scanNumber();
49     Token scanString();
50     Token scanChar();
51     Token scanOperator();
52     TokenType getKeywordType(const std::string& lexeme) const;
53
54 public:
55     Lexer(const std::string& source);
56     Token nextToken();
57     std::vector<Token> tokenize();
58 };
59
60 #endif

```

1.5.2 Implementasi Lexer

Listing 1.4: Implementasi: lexer.cpp (Bagian 1)

```

1 #include "lexer.h"
2 #include <cctype>
3 #include <stdexcept>
4
5 Lexer::Lexer(const std::string& source)
6     : input(source), position(0), line(1), column(1) {
7     // Initialize keywords
8     keywords = {"int", "float", "if", "else",
9                "while", "for", "return"};
10 }
11
12 char Lexer::peek() const {
13     if (position >= input.length()) {
14         return '\0';
15     }
16     return input[position];
17 }
18
19 char Lexer::get() {

```

```

20     if (position >= input.length()) {
21         return '\0';
22     }
23     char c = input[position++];
24     if (c == '\n') {
25         line++;
26         column = 1;
27     } else {
28         column++;
29     }
30     return c;
31 }

```

1.5.3 Handling Whitespace dan Komentar

Listing 1.5: Implementasi: lexer.cpp (Bagian 2 - Whitespace dan Comments)

```

1 void Lexer::skipWhitespace() {
2     while (position < input.length()) {
3         char c = peek();
4         if (std::isspace(c)) {
5             get();
6         } else if (c == '/' && position + 1 < input.length()
7                     && input[position + 1] == '/') {
8             skipLineComment();
9         } else if (c == '/' && position + 1 < input.length()
10                    && input[position + 1] == '*') {
11             skipBlockComment();
12         } else {
13             break;
14         }
15     }
16 }
17
18 void Lexer::skipLineComment() {
19     // Skip "//"
20     get(); get();
21     // Skip until newline or EOF
22     while (peek() != '\n' && peek() != '\0') {
23         get();
24     }
25 }
26
27 void Lexer::skipBlockComment() {
28     // Skip "/*"
29     get(); get();
30     while (position < input.length()) {
31         if (peek() == '*' && position + 1 < input.length()
32             && input[position + 1] == '/') {
33             get(); get(); // Skip "*/"
34             return;
35         }
36         get();

```

```

37     }
38     // Error: unclosed comment
39     throw std::runtime_error("Unclosed block comment at line "
40                             + std::to_string(line));
41 }

```

1.5.4 Scanning Identifier dan Keyword

Listing 1.6: Implementasi: lexer.cpp (Bagian 3 - Identifier)

```

1 Token Lexer::scanIdentifier() {
2     int startLine = line;
3     int startCol = column;
4     std::string lexeme;
5
6     // First character must be letter or underscore
7     if (std::isalpha(peek()) || peek() == '_') {
8         lexeme += get();
9     }
10
11    // Subsequent characters can be alphanumeric or underscore
12    while (std::isalnum(peek()) || peek() == '_') {
13        lexeme += get();
14    }
15
16    // Check if it's a keyword
17    TokenType type = getKeywordType(lexeme);
18    if (type != TokenType::IDENTIFIER) {
19        return Token(type, lexeme, startLine, startCol);
20    }
21
22    return Token(TokenType::IDENTIFIER, lexeme, startLine, startCol);
23 }
24
25 TokenType Lexer::getKeywordType(const std::string& lexeme) const {
26     if (lexeme == "int") return TokenType::KEYWORD_INT;
27     if (lexeme == "float") return TokenType::KEYWORD_FLOAT;
28     if (lexeme == "if") return TokenType::KEYWORD_IF;
29     if (lexeme == "else") return TokenType::KEYWORD_ELSE;
30     if (lexeme == "while") return TokenType::KEYWORD_WHILE;
31     if (lexeme == "for") return TokenType::KEYWORD_FOR;
32     if (lexeme == "return") return TokenType::KEYWORD_RETURN;
33     return TokenType::IDENTIFIER;
34 }

```

1.5.5 Scanning Number Literals

Listing 1.7: Implementasi: lexer.cpp (Bagian 4 - Numbers)

```

1 Token Lexer::scanNumber() {
2     int startLine = line;
3     int startCol = column;

```

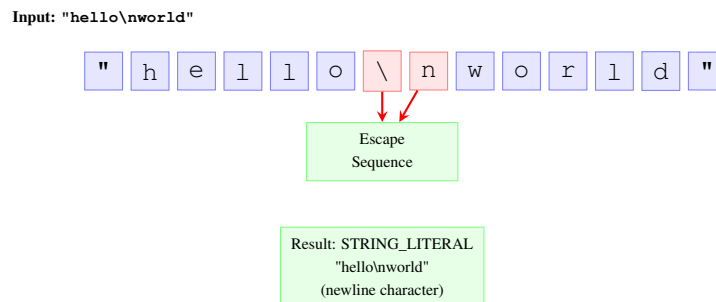
```

4      std::string lexeme;
5      bool isFloat = false;
6
7      // Read integer part
8      while (std::isdigit(peek())) {
9          lexeme += get();
10     }
11
12     // Check for decimal point
13     if (peek() == '.') {
14         lexeme += get();
15         isFloat = true;
16
17         // Read fractional part
18         while (std::isdigit(peek())) {
19             lexeme += get();
20         }
21     }
22
23     // Check for exponent (optional, for future enhancement)
24     if (peek() == 'e' || peek() == 'E') {
25         lexeme += get();
26         if (peek() == '+' || peek() == '-') {
27             lexeme += get();
28         }
29         while (std::isdigit(peek())) {
30             lexeme += get();
31         }
32         isFloat = true;
33     }
34
35     TokenType type = isFloat ? TokenType::FLOAT_LITERAL
36                       : TokenType::INTEGER_LITERAL;
37     return Token(type, lexeme, startLine, startCol);
38 }

```

1.5.6 Scanning String dan Character Literals

Gambar 1.8 menunjukkan contoh handling escape sequences dalam string literal.



Gambar 1.8: Handling escape sequences dalam string literal

Listing 1.8: Implementasi: lexer.cpp (Bagian 5 - Strings)

```

1 Token Lexer::scanString() {
2     int startLine = line;
3     int startCol = column;
4     std::string lexeme;
5
6     // Consume opening quote
7     get(); // Skip opening "
8
9     while (peek() != '"' && peek() != '\0') {
10        if (peek() == '\\') {
11            // Handle escape sequences
12            get(); // Skip backslash
13            char escaped = get();
14            switch (escaped) {
15                case 'n': lexeme += '\n'; break;
16                case 't': lexeme += '\t'; break;
17                case 'r': lexeme += '\r'; break;
18                case '\\': lexeme += '\\'; break;
19                case '"': lexeme += '"'; break;
20                default: lexeme += '\\'; lexeme += escaped; break;
21            }
22        } else {
23            lexeme += get();
24        }
25    }
26
27    if (peek() == '\0') {
28        // Unclosed string
29        return Token(TokenType::INVALID, lexeme, startLine, startCol);
30    }
31
32    get(); // Consume closing "
33    return Token(TokenType::STRING_LITERAL, lexeme, startLine, startCol);
34 }
35
36 Token Lexer::scanChar() {
37     int startLine = line;
38     int startCol = column;
39     std::string lexeme;
40
41     get(); // Skip opening '
42
43     if (peek() == '\\') {
44         // Escape sequence
45         get(); // Skip backslash
46         lexeme += get();
47     } else {
48         lexeme += get();
49     }
50
51     if (peek() != '\\') {
52         return Token(TokenType::INVALID, lexeme, startLine, startCol);
53     }

```

```

54
55     get(); // Consume closing '
56     return Token(TokenType::CHAR_LITERAL, lexeme, startLine, startCol);
57 }

```

1.5.7 Scanning Operators

Listing 1.9: Implementasi: lexer.cpp (Bagian 6 - Operators)

```

1 Token Lexer::scanOperator() {
2     int startLine = line;
3     int startCol = column;
4     char first = get();
5     std::string lexeme(1, first);
6
7     // Check for multi-character operators
8     char next = peek();
9
10    switch (first) {
11        case '=':
12            if (next == '=') {
13                lexeme += get();
14                return Token(TokenType::OP_EQUAL, lexeme, startLine,
15                ↪ startCol);
16            }
17            return Token(TokenType::OP_ASSIGN, lexeme, startLine,
18            ↪ startCol);
19
20        case '!':
21            if (next == '=') {
22                lexeme += get();
23                return Token(TokenType::OP_NOT_EQUAL, lexeme, startLine,
24                ↪ startCol);
25            }
26            return Token(TokenType::OP_NOT, lexeme, startLine, startCol);
27
28        case '<':
29            if (next == '=') {
30                lexeme += get();
31                return Token(TokenType::OP_LESS_EQUAL, lexeme, startLine,
32                ↪ startCol);
33            }
34            return Token(TokenType::OP_LESS, lexeme, startLine, startCol)
35            ↪ ;
36
37        case '>':
38            if (next == '=') {
39                lexeme += get();
40                return Token(TokenType::OP_GREATER_EQUAL, lexeme,
41                ↪ startLine, startCol);
42            }
43            return Token(TokenType::OP_GREATER, lexeme, startLine,
44            ↪ startCol);

```

```

38
39     case '&':
40         if (next == '&') {
41             lexeme += get();
42             return Token(TokenType::OP_AND, lexeme, startLine,
↪ startCol);
43         }
44         return Token(TokenType::INVALID, lexeme, startLine, startCol)
↪ ;
45
46     case '|':
47         if (next == '|') {
48             lexeme += get();
49             return Token(TokenType::OP_OR, lexeme, startLine,
↪ startCol);
50         }
51         return Token(TokenType::INVALID, lexeme, startLine, startCol)
↪ ;
52
53     case '+':
54         return Token(TokenType::OP_PLUS, lexeme, startLine, startCol)
↪ ;
55     case '-':
56         return Token(TokenType::OP_MINUS, lexeme, startLine, startCol)
↪ );
57     case '*':
58         return Token(TokenType::OP_MULTIPLY, lexeme, startLine,
↪ startCol);
59     case '/':
60         return Token(TokenType::OP_DIVIDE, lexeme, startLine,
↪ startCol);
61
62     default:
63         return Token(TokenType::INVALID, lexeme, startLine, startCol)
↪ ;
64     }
65 }

```

1.5.8 Main Tokenization Function

Listing 1.10: Implementasi: lexer.cpp (Bagian 7 - Main Function)

```

1 Token Lexer::nextToken() {
2     skipWhitespace();
3
4     if (position >= input.length()) {
5         return Token(TokenType::END_OF_FILE, "", line, column);
6     }
7
8     char c = peek();
9
10    // Identifier or keyword
11    if (std::isalpha(c) || c == '_') {

```

```
12         return scanIdentifier();
13     }
14
15     // Number
16     if (std::isdigit(c)) {
17         return scanNumber();
18     }
19
20     // String literal
21     if (c == '"') {
22         return scanString();
23     }
24
25     // Character literal
26     if (c == '\\') {
27         return scanChar();
28     }
29
30     // Operators and punctuation
31     if (c == '+' || c == '-' || c == '*' || c == '/' ||
32         c == '=' || c == '!' || c == '<' || c == '>' ||
33         c == '&' || c == '|') {
34         return scanOperator();
35     }
36
37     // Punctuation
38     if (c == ';') {
39         get();
40         return Token(TokenType::SEMICOLON, ";", line, column - 1);
41     }
42     if (c == ',') {
43         get();
44         return Token(TokenType::COMMA, ",", line, column - 1);
45     }
46     if (c == '.') {
47         get();
48         return Token(TokenType::DOT, ".", line, column - 1);
49     }
50     if (c == '(') {
51         get();
52         return Token(TokenType::LPAREN, "(", line, column - 1);
53     }
54     if (c == ')') {
55         get();
56         return Token(TokenType::RPAREN, ")", line, column - 1);
57     }
58     if (c == '{') {
59         get();
60         return Token(TokenType::LBRACE, "{", line, column - 1);
61     }
62     if (c == '}') {
63         get();
64         return Token(TokenType::RBRACE, "}", line, column - 1);
65     }
```

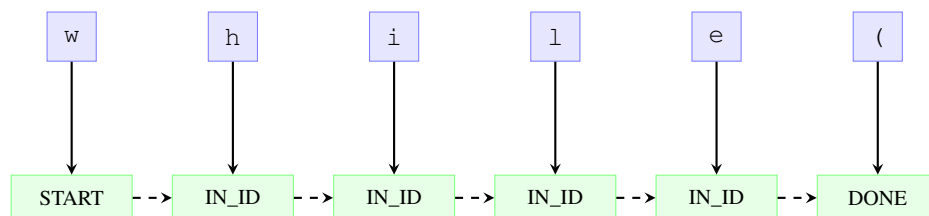


```

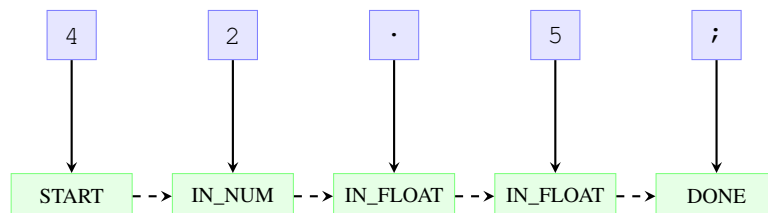
66     if (c == '[') {
67         get();
68         return Token(TokenType::LBRACKET, "[", line, column - 1);
69     }
70     if (c == ']') {
71         get();
72         return Token(TokenType::RBRACKET, "]", line, column - 1);
73     }
74
75     // Unknown character
76     get();
77     return Token(TokenType::INVALID, std::string(1, c), line, column - 1)
78     ↪ ;
79 }
80 std::vector<Token> Lexer::tokenize() {
81     std::vector<Token> tokens;
82     Token token = nextToken();
83     while (token.type != TokenType::END_OF_FILE) {
84         tokens.push_back(token);
85         token = nextToken();
86     }
87     tokens.push_back(token); // Add EOF token
88     return tokens;
89 }

```

Gambar 1.9 menunjukkan proses scanning untuk identifier dan keyword.



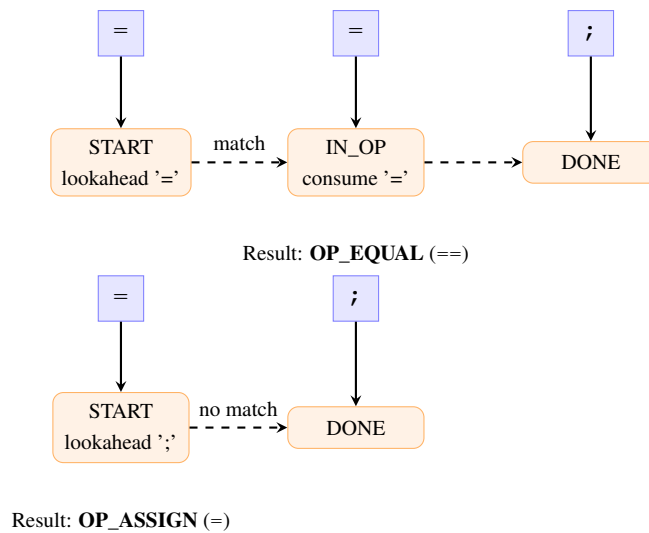
Result: **KEYWORD_WHILE**



Result: **FLOAT_LITERAL 42.5**

Gambar 1.9: Proses scanning keyword/identifier dan literal numerik pada lexer

Gambar 1.10 menunjukkan proses scanning untuk operator multi-character.



Gambar 1.10: Proses scanning operator: perbandingan == dan =

1.6 Error Handling

Error handling dalam lexer harus menangani berbagai kasus edge case:

1.6.1 Unclosed Strings dan Comments

- **Unclosed String:** Jika string literal tidak ditutup sebelum EOF, lexer harus mengembalikan token **INVALID** dengan informasi posisi yang tepat.
- **Unclosed Block Comment:** Jika komentar blok tidak ditutup, dapat di-handle dengan exception atau mengembalikan error token.

1.6.2 Invalid Characters

Karakter yang tidak valid (tidak termasuk dalam kategori token manapun) harus dikembalikan sebagai token **INVALID** dengan informasi posisi untuk error reporting yang baik.

Tabel 1.1 menunjukkan contoh-contoh token yang valid dan tidak valid.

| Input | Token Type | Keterangan |
|------------|-----------------|---------------------------|
| int | KEYWORD_INT | Keyword valid |
| hello | IDENTIFIER | Identifier valid |
| 42 | INTEGER_LITERAL | Integer valid |
| 3.14 | FLOAT_LITERAL | Float valid |
| "hello" | STRING_LITERAL | String valid |
| == | OP_EQUAL | Operator multi-character |
| = | OP_ASSIGN | Operator single-character |
| | INVALID | Karakter tidak valid |
| "unclosed | INVALID | String tidak tertutup |
| /* comment | Error | Comment tidak tertutup |

Tabel 1.1: Contoh token valid dan tidak valid

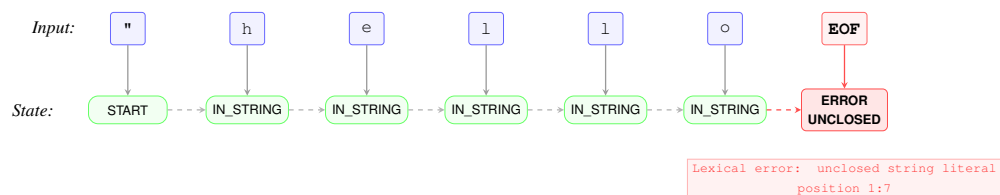
1.6.3 Malformed Numbers

Contoh kasus malformed:

- 123. (titik tanpa digit setelahnya)
- .456 (titik tanpa digit sebelumnya) - dapat di-handle sebagai valid float
- 12.34.56 (multiple decimal points)

Implementasi dapat memilih untuk menerima atau menolak format tertentu sesuai kebutuhan.

Gambar 1.11 menunjukkan contoh error handling untuk unclosed string.



Gambar 1.11: Transisi state lexer dan penanganan kesalahan pada string literal yang tidak tertutup

1.7 Testing Lexer

Unit testing sangat penting untuk memastikan lexer bekerja dengan benar. Berikut contoh test cases:

1.7.1 Test Cases untuk Identifier dan Keyword

Listing 1.11: Test Cases: Identifiers dan Keywords

```
1 void testIdentifiers() {  
2     Lexer lexer("int x = 42;");  
3     Token t1 = lexer.nextToken(); // Should be KEYWORD_INT  
4     Token t2 = lexer.nextToken(); // Should be IDENTIFIER "x"  
5     Token t3 = lexer.nextToken(); // Should be OP_ASSIGN  
6     // ...  
7 }
```

1.7.2 Test Cases untuk Numbers

- 42 → INTEGER_LITERAL
- 3.14 → FLOAT_LITERAL
- 123.456 → FLOAT_LITERAL
- 0 → INTEGER_LITERAL

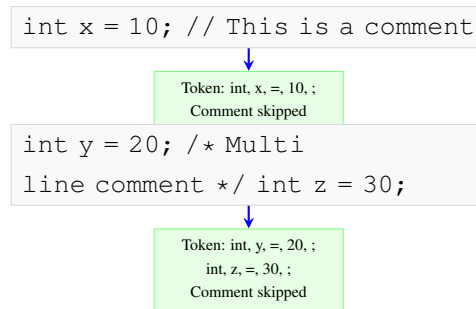
1.7.3 Test Cases untuk Strings

- "hello" → STRING_LITERAL dengan value "hello"
- "hello\nworld" → STRING_LITERAL dengan escape sequence
- "unclosed → INVALID (unclosed string)

1.7.4 Test Cases untuk Comments

- // single line comment → Di-skip, tidak menghasilkan token
- /* multi-line comment */ → Di-skip
- /* unclosed comment → Error atau exception

Gambar 1.12 menunjukkan proses handling komentar dalam lexer.



Gambar 1.12: Handling komentar dalam lexer

1.8 Contoh Penggunaan

Berikut contoh penggunaan lexer untuk tokenize source code sederhana:

Listing 1.12: Contoh Penggunaan Lexer

```

1 #include "lexer.h"
2 #include <iostream>
3
4 int main() {
5     std::string source = R"(
6         int x = 42;
7         float y = 3.14;
8         if (x > 10) {
9             return y;
10        }
11    )";
12
13    Lexer lexer(source);
14    std::vector<Token> tokens = lexer.tokenize();
15
16    for (const auto& token : tokens) {
17        std::cout << "Token: " << token.lexeme
18                  << " Type: " << static_cast<int>(token.type)
19                  << " Line: " << token.line
20                  << " Column: " << token.column << std::endl;
21    }
22
23    return 0;
24 }

```

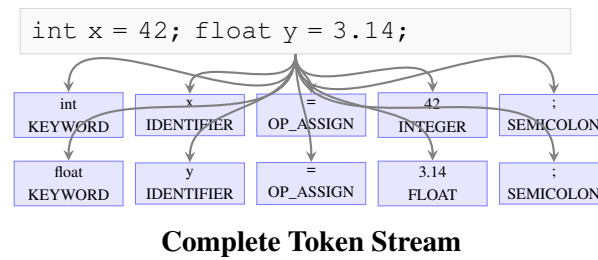
Output yang diharapkan:

```

Token: int Type: 1 Line: 2 Column: 9
Token: x Type: 0 Line: 2 Column: 13
Token: = Type: 13 Line: 2 Column: 15
Token: 42 Type: 8 Line: 2 Column: 17
Token: ; Type: 20 Line: 2 Column: 19
...

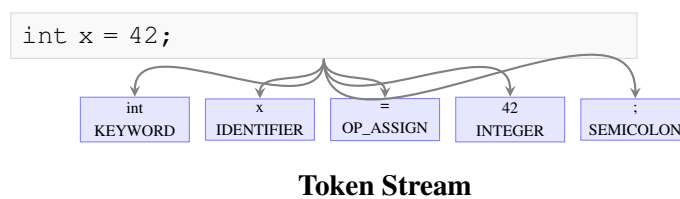
```

Gambar 1.13 menunjukkan contoh lengkap tokenization untuk program sederhana.



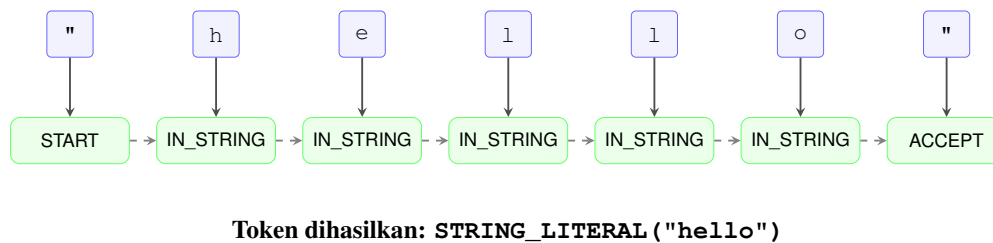
Gambar 1.13: Contoh lengkap tokenization untuk program sederhana

Gambar 1.14 menunjukkan visualisasi token stream untuk contoh kode sederhana.



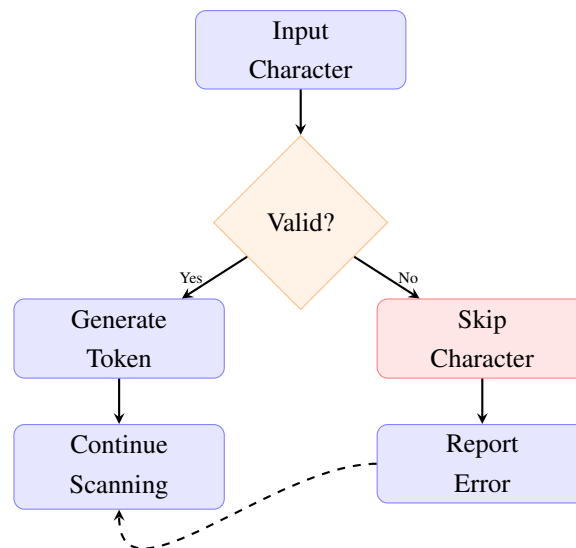
Gambar 1.14: Contoh tokenization: `int x = 42;` menjadi token stream

Gambar 1.15 menunjukkan proses scanning secara detail untuk string "hello".



Gambar 1.15: Representasi formal proses analisis leksikal pada string literal "hello"

Gambar 1.16 menunjukkan strategi error recovery dalam lexer.



Gambar 1.16: Strategi error recovery dalam lexer

1.9 Best Practices

Beberapa best practices dalam implementasi hand-written lexer:

1. **Separation of Concerns:** Pisahkan logika untuk setiap jenis token ke fungsi terpisah
2. **Position Tracking:** Selalu track line dan column untuk error reporting yang baik
3. **Lookahead:** Gunakan `peek()` untuk lookahead tanpa mengkonsumsi karakter
4. **Error Recovery:** Rancang strategi error recovery (misalnya skip invalid character dan lanjut)
5. **Testing:** Buat comprehensive test suite untuk semua edge cases
6. **Documentation:** Dokumentasikan token types dan format yang didukung

1.10 Kesimpulan

Dalam bab ini, kita telah mempelajari:

1. Struktur token dan token types untuk subset bahasa C
2. Konsep finite state machine dalam konteks lexical analysis
3. Implementasi hand-written lexer dalam C++ dengan handling:

- Identifier dan keyword recognition
 - Number literals (integer dan float)
 - String dan character literals dengan escape sequences
 - Operators (single dan multi-character)
 - Whitespace dan komentar (single-line dan multi-line)
4. Error handling untuk edge cases
 5. Testing strategies untuk lexer

Implementasi hand-written lexer memberikan pemahaman mendalam tentang proses tokenization dan menjadi dasar untuk memahami bagaimana lexer generator seperti Flex bekerja di belakang layar.

1.11 Referensi dan Bahan Bacaan Lanjutan

Untuk memperdalam pemahaman tentang implementasi lexer, mahasiswa disarankan membaca:

- **Dragon Book:** Aho, Lam, Sethi, & Ullman (2006). *Compilers: Principles, Techniques, and Tools* ? - Bab 3: Lexical Analysis
- **Engineering a Compiler:** Cooper & Torczon (2011) ? - Bab 2: Scanning
- **OpenGenus - Build Lexer:** Tutorial tentang hand-written lexer ?
- **Aoyama Gakuin University:** Lecture notes tentang lexical analysis ?
- **GeeksforGeeks:** Contoh implementasi lexical analyzer dalam C++ ¹
- **Programming Notes:** Tutorial tentang simple lexer menggunakan finite state machine ²

¹<https://www.geeksforgeeks.org/cpp/lexical-analyzer-in-cpp/>

²<https://www.programmingnotes.org/4699/cpp-simple-lexer-using-a-finite-state-machine/>