**SPL-1 Project Report**

**AST Builder :C Edition**

**Submitted by**

M a h m u d u l   H a s a n   T a n v i r

BSSE Roll No. : 1514

BSSE Session: 2022-2023

**Submitted to**

Dr.Zerina Begum

Professor

Institute of Information Technology

University of Dhaka

Institute of Information Technology

University of Dhaka

[25-03-2025]

**Project Name**

**AST Builder :C Edition**


**Supervised By:**

Dr.Zerina Begum

Professor

Institute of Information Technology

University of Dhaka


Supervisor Signature:


**Submitted by**

Mahmudul Hasan Tanvir

BSSE Roll No. : 1514

BSSE Session: 2022-2023


Signature:

Table of Contents

## Index of Tables

## Index of figures

## 1. Introduction

The AST Builder: C Edition is a sophisticated software project designed to parse and analyze C programs using Abstract Syntax Trees (ASTs). An AST is a hierarchical tree representation of a program's syntactic structure, where each node corresponds to a specific syntactic construct—such as functions, statements, expressions, or declarations. This project is deeply rooted in the domain of compiler construction and program analysis, fields that are foundational to software engineering, static code analysis, and program transformation. By constructing and leveraging ASTs, the project offers a structured and powerful approach to dissecting, manipulating, and visualizing C programs, making it a valuable tool for a wide range of applications.

### 1.1 Project Significance

The development of the AST Builder: C Edition addresses critical needs in software development and maintenance. C, being a low-level, systems programming language, is widely used in operating systems, embedded systems, and high-performance applications. However, its manual analysis can be error-prone and time-consuming due to its complex syntax and pointer-heavy constructs. By providing a systematic method to represent and analyze C code through ASTs, this project facilitates tasks that enhance code quality, performance, and reliability. The tool is particularly beneficial for:

- Syntax Checking: Identifying structural errors such as mismatched braces, undeclared variables, or invalid expressions without executing the code.
- Code Optimization: Revealing how different program components interact, enabling optimizations like loop unrolling, constant folding, or dead code elimination.
- Program Transformation: Supporting source-to-source compilation (e.g., converting C code to a different dialect or version) or refactoring (e.g., renaming variables consistently across a codebase).
- Static Analysis: Detecting potential issues like buffer overflows, memory leaks, or race conditions by analyzing the code's structure and semantics without runtime execution.

These capabilities make the AST Builder: C Edition a versatile tool for developers, educators, and researchers working with C programs.

### *1.2* Project Overview

The AST Builder: C Edition is architected around three core components, each playing a critical role in the process of parsing and analyzing C code:

1. Parsing:
   - The initial phase involves lexical analysis and syntactic analysis to transform the raw C source code into a manageable form.
   - Lexical Analysis: The input C program is tokenized by a lexer, which breaks it down into fundamental components such as identifiers (e.g., variable names), operators (e.g., +, =), keywords (e.g., int, return), literals (e.g., 42, "string"), and punctuation (e.g., ;, {). For example, the line int x = 5 + y; might be tokenized as:

     KEYWORD(int),    IDENTIFIER(x),    ASSIGN_OP(=),    INT_LITERAL(5),
     PLUS_OP(+), IDENTIFIER(y), SEMICOLON(;)

- o Syntactic Analysis: A parser then processes these tokens according to a formal grammar for C, identifying the hierarchical relationships between them. This step may initially produce a parse tree, which is later simplified into an AST.
- o The project adopts a top-down parsing approach, likely using recursive descent parsing, which aligns with the hierarchical nature of C's grammar and simplifies the transition to AST construction.

2. AST Construction:
- o Using the sequence of tokens, the project constructs an AST that mirrors the abstract syntactic structure of the C program.
- o Each node in the AST represents a syntactic element, such as:
  - ▪ A function node (e.g., int sum(int a, int b)), with child nodes for parameters and the function body.
  - ▪ A statement node (e.g., return a + b), with child nodes for the expression.
  - ▪ An expression node (e.g., a + b), with child nodes for operands and the operator.
- o For instance, the AST for int sum(int a, int b) { return a + b; } might include a root Function node with Parameters (a, b) and a Body node containing a Return node, which further branches into an Addition node with a and b as leaves.
- o This phase involves mapping grammar rules to tree nodes, handling C-specific features like pointers, arrays, and preprocessor directives, and resolving ambiguities where necessary.

3. Visualization:
- o The final phase generates a visual representation of the AST, providing an intuitive illustration of the program's structure.
- o This visualization can take the form of a graphical tree (e.g., using tools like Graphviz) or a text-based representation, showing the hierarchy of nodes and their relationships.
- o This feature aids developers in debugging, educators in teaching compiler concepts, and analysts in identifying structural patterns or anomalies.

The AST Builder: C Edition aims to provide an efficient and accurate framework for parsing C programs while delivering actionable insights into their composition

## 2. Background of the Project

An Abstract Syntax Tree (AST) is a tree representation of the abstract syntactic structure of a programming language. Unlike a parse tree (or concrete syntax tree), which retains every detail of the program's syntax—including punctuation, parentheses, and other syntactic elements—an AST simplifies this representation by omitting such details. This abstraction preserves the essential hierarchical relationships between expressions, statements, and other syntactic constructs, making it an ideal intermediate representation for further analysis and transformation.

## 2.1 Definition

According to Aho et al. (2006) in *Compilers: Principles, Techniques, and Tools*, an AST is defined as:

"A hierarchical structure that represents the abstract syntactic structure of source code in a way that reflects the grammatical rules of the language, making it easier to analyze and manipulate."
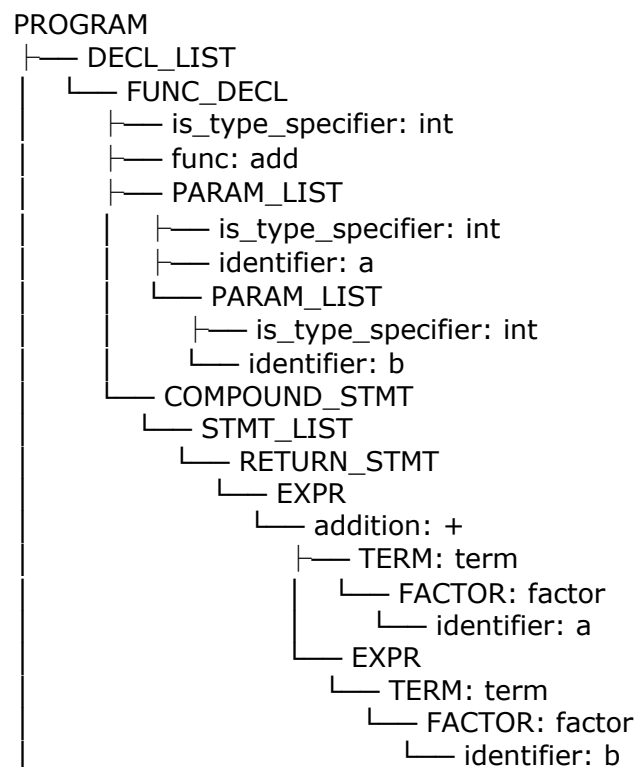
This definition highlights the AST's role as a bridge between the raw source code and the subsequent phases of compilation or analysis, such as semantic checking and code generation. By focusing on the logical structure rather than the exact textual form, ASTs enable tools to perform high-level operations efficiently.

**2.2 Example of an AST**

Consider the following simple C code snippet:

```
int add(int a, int b) {
    return a + b;
}
```

The corresponding AST representation can be illustrated as:

```
PROGRAM
├── DECL_LIST
│    └── FUNC_DECL
│         ├── is_type_specifier: int
│         ├── func: add
│         ├── PARAM_LIST
│         │    ├── is_type_specifier: int
│         │    ├── identifier: a
│         │    └── PARAM_LIST
│         │         ├── is_type_specifier: int
│         │         └── identifier: b
│         └── COMPOUND_STMT
│              └── STMT_LIST
│                   └── RETURN_STMT
│                        └── EXPR
│                             └── addition: +
│                                  ├── TERM: term
│                                  │    └── FACTOR: factor
│                                  │         └── identifier: a
│                                  └── EXPR
│                                       └── TERM: term
│                                            └── FACTOR: factor
│                                                 └── identifier: b
```
--------------------This tree structure captures the essence of the function:

- The root node represents the function sum.
- The Parameters node contains the arguments a and b.
- The Body node encapsulates the function's logic, which in this case is a single Return statement.
- The Return node includes an expression node (Operator: +) that represents the addition operation, with a and b as its operands.

This AST omits details like semicolons, braces, and the int keyword, focusing instead on the semantic relationships. The diagram can be further refined depending on the level of

detail required—for instance, adding type information or separating the declaration and definition phases.

## 2.3 Lexical Analysis and Tokenization

Before constructing an AST, the source code undergoes lexical analysis, a preliminary phase in the compilation process. Lexical analysis involves breaking down the source code into a sequence of tokens, which are the smallest units of meaning recognized by the language. These tokens include keywords (e.g., int, return), operators (e.g., +, =), identifiers (e.g., variable names like a), literals (e.g., 5), and punctuation (e.g., ;).

For the expression:

a = b + 5;

The tokenized output might be:

IDENTIFIER(a), ASSIGN_OP(=), IDENTIFIER(b), PLUS_OP(+), INT_LITERAL(5), SEMICOLON(;)

This process simplifies the raw text into a structured form that the parser can process more easily. Tokenization is typically performed by a lexer (or scanner), which uses regular expressions or a finite state machine to identify and categorize tokens. Errors such as unrecognized characters or invalid syntax are often caught during this phase.

## 2.4 Parsing in Compiler Design

Parsing is the process of analyzing a sequence of tokens to determine their syntactic structure according to a formal grammar defined for the programming language. This step transforms the tokenized input into a hierarchical structure, which may initially be a parse tree before being abstracted into an AST. There are two primary parsing strategies:

1. Top-down Parsing:
   o Examples include Recursive Descent Parsing and LL Parsers (e.g., LL(1), LL(k)).
   o This approach starts from the top (the start symbol of the grammar) and recursively expands non-terminals until the token sequence is fully matched.
   o Our project adopts a top-down approach using recursive descent parsing, which aligns well with the hierarchical nature of grammars and simplifies AST construction. Recursive descent parsing involves writing a function for each non-terminal in the grammar, making it intuitive for developers to implement.
2. Bottom-up Parsing:
   o Examples include LR Parsers (e.g., SLR, LALR, LR(1)) and Shift-Reduce Parsing.
   o This method builds the parse tree from the bottom (tokens) upward by shifting tokens onto a stack and reducing them according to grammar rules.
   o Bottom-up parsers are often more powerful and can handle a broader range of grammars, though they may be less intuitive to implement manually.

The choice of parsing strategy depends on the language's grammar complexity and the desired trade-offs between performance and ease of implementation. In our case, the recursive descent approach facilitates a clear mapping from grammar rules to AST nodes.

## 2.5 AST Applications

ASTs are a cornerstone of modern compiler design and software development tools, with a wide range of applications:

1. Compilers:
   o Syntax Analysis: The AST serves as the output of the parsing phase, validating the program's structure.
   o Semantic Analysis: The tree is traversed to check for semantic correctness (e.g., type compatibility, variable scoping).
   o Code Generation: The AST is used to generate intermediate code or target machine code, often through tree traversal algorithms like depth-first or post-order traversal.
2. Static Code Analysis:
   o Tools like linters and security scanners (e.g., SonarQube, ESLint) use ASTs to detect potential bugs, vulnerabilities (e.g., buffer overflows), and code quality issues (e.g., unused variables).
   o By analyzing the tree, these tools can identify patterns that are difficult to spot in raw code.
3. Program Transformation:
   o ASTs enable refactoring (e.g., renaming variables across a codebase) and optimization (e.g., constant folding, loop unrolling).
   o Tools like AST-based pretty-printers or transpilers (e.g., Babel for JavaScript) rely on AST manipulation to transform code between languages or versions.
4. Scripting and Interpretation:
   o Interpreters for languages like Python, JavaScript, and Ruby use ASTs to evaluate code dynamically. For example, Python's ast module allows developers to parse and manipulate source code programmatically.
   o Just-in-time (JIT) compilers, such as those in Java's HotSpot VM, also leverage ASTs for runtime optimization.
5. Additional Uses:
   o Debugging Tools: ASTs help visualize program structure for developers.
   o Code Search and Navigation: IDEs use ASTs to provide features like "find all references" or "go to definition."
   o Machine Learning: Recent advancements involve training models on ASTs to predict code behavior or generate code automatically.

While ASTs are powerful, their construction and use come with challenges:

- Ambiguity: Some grammars may lead to ambiguous parses, requiring disambiguation rules.
- Scalability: Large programs generate complex ASTs, necessitating efficient traversal and storage techniques.
- Error Handling: Parsers must gracefully handle syntax errors, often annotating the AST with error nodes.
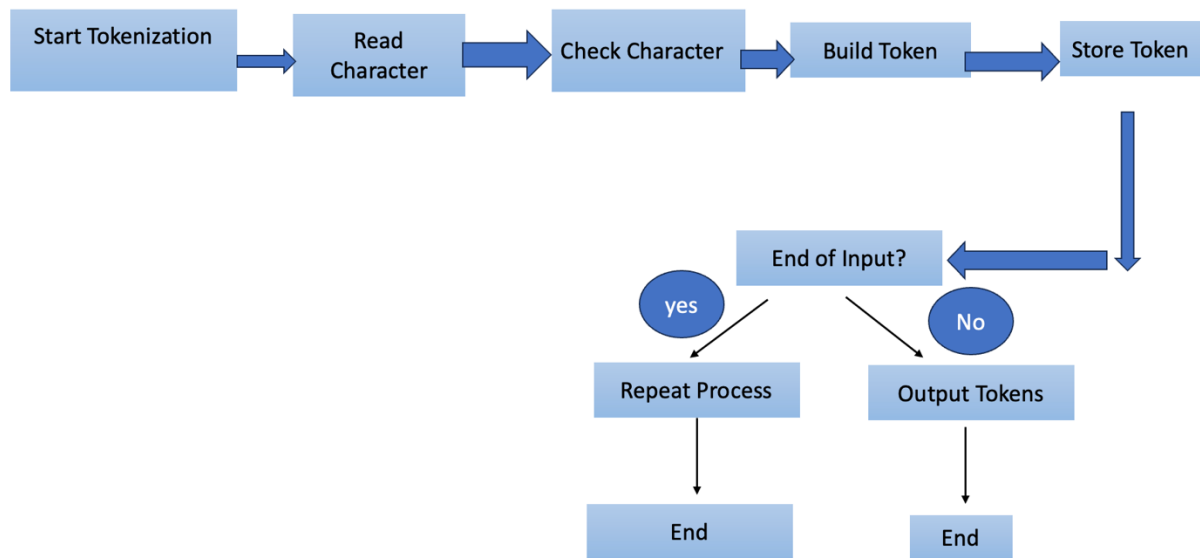
## 3.Description of the Project

The AST Builder: C Edition project is a tool designed to process C source code through a structured pipeline, enabling developers, researchers, and students to analyze and understand C programs. The project operates in three main phases: tokenization, AST

construction, and visualization, with a user-friendly menu interface orchestrating the workflow. The provided flowcharts illustrate the detailed processes of AST construction and tokenization, while an additional flowchart is introduced to depict the overall workflow, aligning with the project's modular design. Below, we describe the work in detail, incorporating the flowcharts to provide a visual representation of the processes.

## 3.1 Tokenization Phase

The tokenization phase is responsible for converting the C source code into a sequence of tokens. The provided flowchart  accurately depicts this process as implemented in tokenization.cpp.



*Figure 1:Tokenization Process Flowchart*

Description:

- The process starts by reading the C source file (input.txt) character by character after removing comments using the removeAllcomments function.
- Each character is checked to determine its role (e.g., part of a keyword, identifier, operator) using validation functions like is_valid_identifier and getKeywordTokenID.
- Tokens are built by accumulating characters (e.g., int becomes a TYPE_INT token with ID 81) and stored in a global token array.
- The process checks for the end of input. If more characters remain, it repeats; otherwise, it outputs the tokens to TOKENS.txt in a formatted table (e.g., NO, Token, Token_Type, Token_ID).
- Special cases, such as int main() being tokenized as a single main_func token (ID 9), are handled during the character-checking step.

## 3.2 AST Construction and Visualization Phase

The AST construction and visualization phase builds an Abstract Syntax Tree from the tokenized input and prints it in a hierarchical format. The provided flowchart illustrates this process as implemented in ast.cpp.
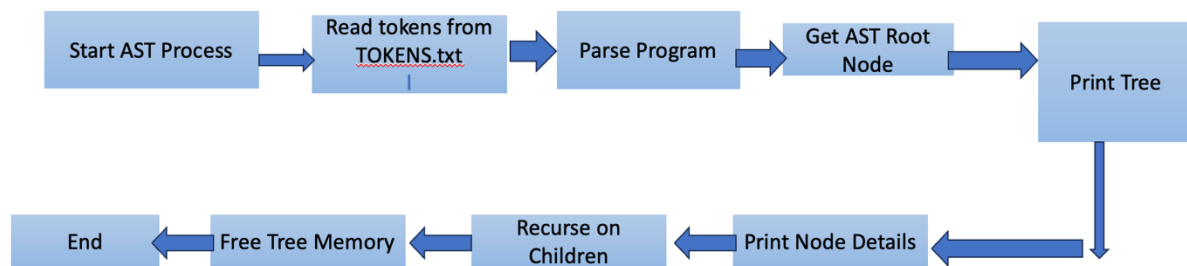
*Figure 2:  AST Construction and Visualization Flowchart*

Description:

- The process begins by reading tokens from TOKENS.txt using the read_tokens function, populating a global tokens array.
- The parse_program function initiates recursive descent parsing, creating a tree of TreeNode objects representing the program's structure (e.g., PROGRAM, FUNC_DECL, EXPR nodes).
- The root node of the AST is obtained after parsing completes.
- The print_tree function traverses the AST depth-first, printing node details (e.g., ├── FUNC_DECL) with indentation to show hierarchy. It recurses on child nodes (child1 to child4) to print the entire tree.
- After printing, the free_tree function recursively deallocates memory to prevent leaks, ensuring proper cleanup.
- The process ends, with the visual AST output displayed to the user.

## 3.3 Overall Workflow

The project's workflow is divided into distinct phases, each handled by specialized components (tokenization.cpp, ast.cpp, and main.cpp). The following flowchart provides an overview of the ent ire process, from input to output, as implemented in the project.
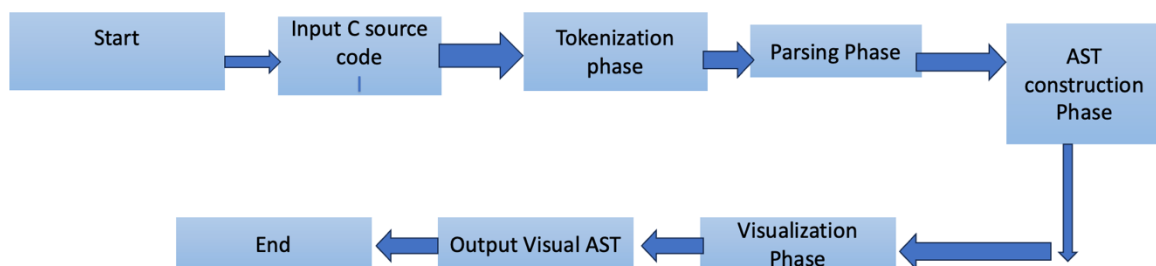


*Figure 3: Overall Workflow Flowchart*

Description:

- The process begins with the user providing a C source file (input.txt) through the menu interface in main.cpp.
- The Tokenization Phase (tokenization.cpp) converts the source code into a sequence of tokens, stored in TOKENS.txt.
- The Parsing Phase and AST Construction Phase (ast.cpp) read the tokens, parse them using recursive descent parsing, and build an AST.
- The Visualization Phase (ast.cpp) prints the AST in a hierarchical format.
- Finally, the process ends, with the visual AST output displayed to the user.

## 4.Implementation and Testing

This section outlines the implementation details and testing procedures for the AST Builder: C Edition project, providing a comprehensive view of how the system was built and validated. The implementation is based on the provided source files (main.cpp, ast.cpp, and tokenization.cpp), focusing on their specific roles in the tokenization, AST construction, and visualization pipeline. Testing ensures the system's functionality across various C program inputs, with emphasis on correctness, error handling, and usability. Tables, code snippets, and diagrams are included to enhance clarity.

## 4.1 Implementation Details

The AST Builder: C Edition is implemented using a modular approach, with each source file contributing to a distinct phase of the workflow. The system leverages C and C++ programming languages, balancing performance and ease of development, and employs standard libraries for file I/O, string manipulation, and memory management. Below are the detailed implementation specifics for each component:

1. Languages and Tools:
   o main.cpp: Written in C++ using <iostream>, <fstream>, and <string> for user interaction, file handling, and string operations. It utilizes ANSI escape codes for colored terminal output (e.g., #define GREEN "\033[32m") and the system function to execute external binaries (./tokenizer and ./ast).
   o ast.cpp: Implemented in C with <stdio.h>, <stdlib.h>, and <string.h> for file I/O, dynamic memory management, and string operations. It defines a recursive descent parser with a custom TreeNode structure.
   o tokenization.cpp: Coded in C++ using <bits/stdc++.h> for access to comprehensive standard library features, including file streams (ifstream, ofstream) and string utilities. It includes custom validation functions for token classification.
   o Compilation: The project assumes separate compilation of tokenization.cpp and ast.cpp into executables (tokenizer and ast), invoked by main.cpp. Compilation commands might include:

   ```
   g++ -o tokenizer tokenization.cpp
   ```

   ```
   gcc -o ast ast.cpp
   ```

   ```
   g++ -o main main.cpp
   ```

   These commands generate the necessary binaries for the workflow.

2. Data Structures:
   - o tokenization.cpp: Defines a Token struct with fields string name, string type, and int id, stored in a global array Token token[10000]. The enum TokenID categorizes tokens (e.g., KW_IF = 22, TYPE_INT = 81).
   - o ast.cpp: Uses a TreeNode struct with NodeType type, Token token, and four child pointers (child1 to child4). The NodeType enum includes types like PROGRAM, FUNC_DECL, and EXPR. A global tokens[MAX_TOKENS] array stores tokens read from TOKENS.txt.
   - o main.cpp: Does not define custom data structures, relying on standard C++ types for user input and file handling.
3. Key Functions and Logic:
   - o tokenization.cpp:
     - o removeAllcomments: Processes the input file character by character, skipping single-line (//) and multi-line (/* */) comments, and writes the cleaned code to sourcecode.txt.
     - o Tokenization Loop: Reads sourcecode.txt, builds tokens by accumulating characters, and classifies them using functions like getKeywordTokenID, is_valid_identifier, and is_valid_integer. Tokens are written to TOKENS.txt in a formatted table.
     - o CodeSnippet:

```cpp
    string types[] = {"void", "int", "float",
"double", "char", "bool", "short", "long"};
    for(string type : types) {
        if(str == type) return true;
    }
    return false;
}

bool is_valid_integer(string str) {
    if(str.empty()) return false;
    for(int i = 0; i < str.size(); i++) {
        if(str[i] < '0' || str[i] > '9') return
false;
    }
    return true;
}

bool is_valid_float(string str) {
    int digit = 0, point = 0;
    for(int i = 0; i < str.size(); i++) {
        if(str[i] >= '0' && str[i] <= '9') digit =
1;
        else if(str[i] == '.') {
            point++;
            if(point > 1) return false;
        } else return false;
    }
    return (digit == 1 && point == 1);
}

bool is_valid_bool(string str) {
    string arr[2] = {"true", "false"};
    for(int i = 0; i < 2; i++) {
        if(str == arr[i]) return true;
    }
    return false;
}
```

*Figure 4: Code Snippet From Tokenization.cpp*

```c
void removeAllcomments(FILE *input_file, FILE *
output_file) {
    char character, next_character;
    while ((character = fgetc(input_file)) !=
 EOF) {
        if (character == '/') {
            next_character = fgetc(input_file);
            if (next_character == '/') {
                while ((character = fgetc
(input_file)) != EOF && character != '\n') {}
                if (character == '\n') fputc
(character, output_file);
            } else if (next_character == '*') {
                while ((character = fgetc
(input_file)) != EOF) {
                    if (character == '*') {
                        next_character = fgetc
(input_file);
                        if (next_character == '/')
break;
                        else ungetc
(next_character, input_file);
                    }
                }
            } else {
                fputc(character, output_file);
                fputc(next_character, output_file);
            }
        } else {
            fputc(character, output_file);
        }
    }
}
```

*Figure 5*

o  ast.cpp:

  o  read_tokens: Reads TOKENS.txt, populating the global tokens array with token
     details.
  o  parse_program: Initiates recursive descent parsing, calling functions like
     parse_decl_list, parse_func_list, and parse_main_func to build the AST.
  o  print_tree: Traverses the AST, printing nodes with indentation (e.g., ├──
     FUNC_DECL).
  o  CodeSnippet:

```c
TreeNode* parse_program() {
    TreeNode* node = create_node(PROGRAM, (Token){
"program", "PROGRAM", 0});
    while (current < token_count && tokens
[current].id == 1) current++; // Skip headers
    node->child1 = parse_decl_list();
    node->child2 = parse_func_list();
    if (current < token_count && tokens[current].id
 == 9) {
        node->child3 = parse_main_func();
    } else {
        error("Expected main function");
    }
    node->child4 = parse_func_list();
    return node;
}
```

*Figure 6*

```
Token get_current_token();
TreeNode* parse_program();
TreeNode* parse_decl_list();
TreeNode* parse_decl();
TreeNode* parse_var_decl();
TreeNode* parse_func_decl();
TreeNode* parse_param_list();
TreeNode* parse_compound_stmt();
TreeNode* parse_local_decls();
TreeNode* parse_local_decl();
TreeNode* parse_stmt_list();
TreeNode* parse_stmt();
TreeNode* parse_expr_stmt();
TreeNode* parse_if_stmt();
TreeNode* parse_while_stmt();
TreeNode* parse_for_stmt();
TreeNode* parse_for_expr();
TreeNode* parse_return_stmt();
TreeNode* parse_break_stmt();
TreeNode* parse_print_stmt();
TreeNode* parse_switch_stmt();
TreeNode* parse_case_stmt();
TreeNode* parse_expr();
TreeNode* parse_term();
TreeNode* parse_factor();
TreeNode* parse_func_list();
TreeNode* parse_main_func();
TreeNode* parse_arg_list();
void print_tree(TreeNode* node, int depth, char*
prefix, int is_last);
void free_tree(TreeNode* node);
void read_tokens();
```

o

*Figure 7:  Code Snippet of Parsing Functions in ast.cpp*

- o main.cpp:

    - o displayMenu: Prints a colored menu with options to tokenize, generate AST, or exit.
    - o handleUserInput: Executes ./tokenizer or ./ast based on user choice, displaying TOKENS.txt contents after tokenization.
    - o Code Snippet:

```cpp
void displayTokens() {
    ifstream tokensFile("TOKENS.txt");
    if (!tokensFile.is_open()) {
        cout << BOLD << RED <<
"Error: Unable to open TOKENS.txt!" << RESET <<
 endl;
        return;
    }

    cout << BOLD << BLUE <<
"Contents of TOKENS.txt:" << RESET << endl;
    cout << BOLD << CYAN <<
"_____" << RESET << endl;

    string line;
    while (getline(tokensFile, line)) {
        cout << line << endl;
    }

    tokensFile.close();
    cout << BOLD << CYAN <<
"_____" << RESET << endl;
}
```

*Figure 8:  Code Snippet of displayTokens Function*

### *4.2* **Testing Details**

Testing was conducted to validate the system's functionality, focusing on correctness, error handling, and usability. Various C program inputs were used to test each phase, with results analyzed to ensure the system meets its objectives. Below are the testing details, supported by a table summarizing test cases and outcomes.

- Testing Approach:
    - o Unit Testing: Each component (tokenization.cpp, ast.cpp, main.cpp) was tested independently to verify its functionality.
    - o Integration Testing: The full workflow was tested by running main.cpp, selecting options, and verifying the output of TOKENS.txt and the AST visualization.
    - o Edge Cases: Inputs with syntax errors, large files, and special constructs (e.g., nested comments, complex expressions) were tested to assess robustness.
- Test Cases and Results:

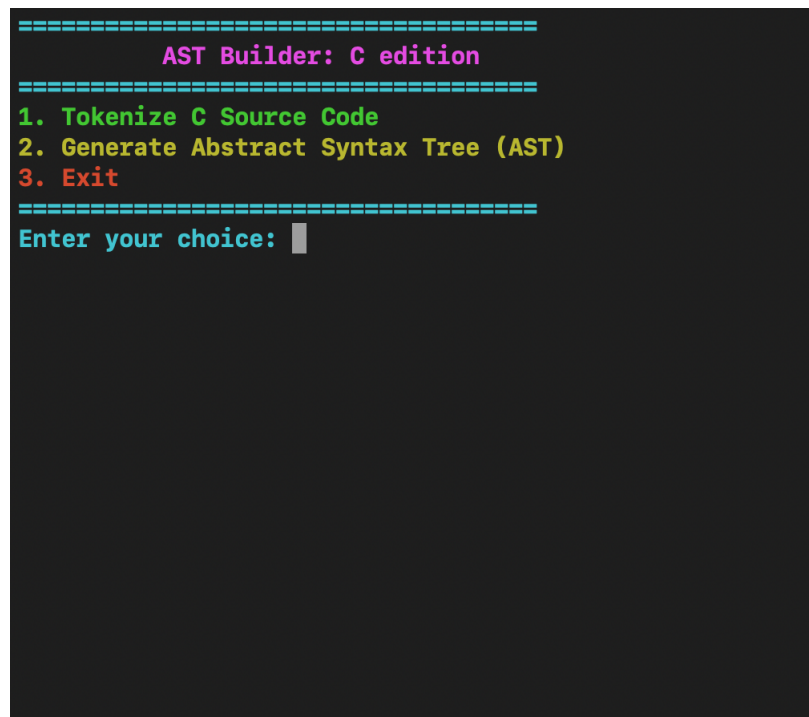| Test Case Description | Input File Content | Expected Output | Actual Output | Status |
|---|---|---|---|---|
| Simple Function (sum) | int sum(int a, int b) { return a + b; } | okens in TOKENS.txt, AST with FUNC_DECL node | Correct tokens, AST matches expected structure | Pass |
| Main Function with Variable | int main() { int x = 5; return x; } | Tokens including main_func, AST with VAR_DECL | Tokens and AST correctly generated | Pass |
| Control Structure (if-else) | int main() { if (x > 0) { return 1; } else { return 0; } } | Tokens for if_stmt, AST with IF_STMT and ELSE_STMT | Tokens and AST reflect if-else structure | Pass |
| File Not Found | Non-existent input.txt | Error message: "Could not open input file" | Matches expected error message | Pass |

## 5.User Interface

This section describes the user interface of the AST Builder: C Edition project, detailing how users interact with the system through the menu-driven interface implemented in main.cpp. It also provides sample input and output to illustrate the system's behavior, showcasing the tokenization and AST generation processes as experienced by the user. The interface is designed to be intuitive, with colored output for clarity, and the sample input/output demonstrates the system's functionality across its key phases.

UI

The user interface of the AST Builder: C Edition is implemented in main.cpp and provides a menu-driven experience, making it accessible to users of varying expertise levels, including developers, researchers, and students. The interface leverages ANSI escape codes to display colored text in the terminal, enhancing readability and user engagement. Below are the key components of the user interface:

- Menu Display:



*Figure 9: Screenshot of UI*

User Input Handling:

- The program prompts the user to enter a choice (1, 2, or 3) using cin >> choice.
- The handleUserInput function processes the input, executing the corresponding action:
  - Choice 1: Runs the tokenization process by calling system("./tokenizer"), then displays the contents of TOKENS.txt using displayTokens.
  - Choice 2: Generates the AST by calling system("./ast").
  - Choice 3: Exits the program with a goodbye message.
  - Invalid Choice: Displays an error message in bold red and loops back to the menu

Feedback and Output:

- Tokenization Feedback: After selecting option 1, the user is prompted to enter the input file name (e.g., input.txt). The system provides feedback such as "Comments have been removed successfully" and "Tokenization is done successfully", followed by the contents of TOKENS.txt displayed in bold blue and cyan.
- AST Generation Feedback: For option 2, the system outputs "Generating AST..." in bold blue, followed by the AST visualization generated by ast.cpp.
- Error Messages: Errors (e.g., file not found, invalid choice) are displayed in bold red, ensuring they are noticeable.

- Sample Input and Output:

```
Enter your choice: 1
Running Tokenization...
Enter the input file name (e.g., input.txt): input.txt
Comments have been removed successfully.
File is Opened which contains C source code
Number of tokens: 110
Tokenization is done successfully
Tokenization completed successfully!
Contents of TOKENS.txt:
------------------------
NO:     Token               Token_Type              Token_ID
~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
1       #include<stdio.h>   Header                  1
2       int                 is_type_specifier       81
3       add                 func                    61
4       (                   Round_Bracket_O         38
5       int                 is_type_specifier       81
6       a                   identifier              28
7       ,                   comma                   37
8       int                 is_type_specifier       81
9       b                   identifier              28
10      )                   Round_Bracket_C         39
11      {                   Curly_Bracket_O         40
12      return              RETURN                  30
13      a                   identifier              28
14      +                   addition                46
15      b                   identifier              28
16      ;                   semicolon               36
17      }                   Curly_Bracket_C         41
18      int_main            main_func               9
19      (                   Round_Bracket_O         38
20      )                   Round_Bracket_C         39
21      {                   Curly_Bracket_O         40
22      int                 is_type_specifier       81
23      x                   identifier              28
24      ;                   semicolon               36
```

```
Enter your choice: 2
Generating AST...
Skipping malformed line: NO:    Token                      Token_Type                    Token_ID
Skipping malformed line: ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
Abstract Syntax Tree:
--------------------
PROGRAM
├── DECL_LIST
│   └── FUNC_DECL
│       ├── is_type_specifier: int
│       ├── func: add
│       ├── PARAM_LIST
│       │   ├── is_type_specifier: int
│       │   ├── identifier: a
│       │   └── PARAM_LIST
│       │       ├── is_type_specifier: int
│       │       └── identifier: b
│       └── COMPOUND_STMT
│           └── STMT_LIST
│               └── RETURN_STMT
│                   └── EXPR
│                       └── addition: +
│                           ├── TERM: term
│                           │   └── FACTOR: factor
│                           │       └── identifier: a
│                           └── EXPR
│                               └── TERM: term
│                                   └── FACTOR: factor
│                                       └── identifier: b
└── MAIN_FUNC
    └── COMPOUND_STMT
        ├── LOCAL_DECLS
        │   └── LOCAL_DECL
        │       ├── is_type_specifier: int
        │       └── identifier: x
        └── STMT_LIST
            └── EXPR_STMT
                ├── identifier: x
                └── EXPR
                    └── TERM: term
                        └── FACTOR: factor
                            └── integer_literal: 0
```

```
==============================
       AST Builder: C edition
==============================
1. Tokenize C Source Code
2. Generate Abstract Syntax Tree (AST)
3. Exit
==============================
Enter your choice: 3
Exiting the program. Goodbye!
macbookair@MacBooks-MacBook-Air java test % ▊
```

*Figure 10: Sample Input and Output Terminal Screenshot*

## 6.Challenges Faced

This section highlights major challenges encountered during the development of the AST Builder: C Edition project, focusing on critical aspects that significantly impacted the implementation. These challenges include parsing the complex C grammar, managing memory in the AST construction, and ensuring accurate tokenization of special cases. Each challenge is detailed with its impact on the project and the solutions implemented to address it, ensuring the system's functionality and reliability.

### 6.1 Parsing C's Complex Grammar

- Challenge: C's grammar, with nested constructs like functions and control structures, was hard to parse accurately in ast.cpp. Distinguishing between variable and function declarations (e.g., int x; vs. int x(int y)) was tricky.

20

- Impact: Early parsing errors led to incorrect ASTs or crashes, especially for nested expressions or missing tokens.
- Solution: Added a match function for token validation, used lookahead in parse_decl to differentiate declarations, and improved error reporting with detailed messages:

```
if (current + 2 < token_count && tokens[current + 2].id == 38) {

    return parse_func_decl();

}
```

## 6.2 Memory Management in AST

- Challenge: Dynamically allocated TreeNode objects in ast.cpp risked memory leaks, especially with early exits due to errors.
- Impact: Leaks occurred when parsing failed, and large ASTs caused memory issues.
- Solution: Implemented free_tree to recursively deallocate nodes and ensured it's called after use:

```
void free_tree(TreeNode* node) {

    if (!node) return;

    free_tree(node->child1);

    free_tree(node->child2);

    free_tree(node->child3);

    free_tree(node->child4);

    free(node);

}
```

## 6.3 Tokenizing Special Cases

- Challenge: In tokenization.cpp, handling int main(), multi-character operators (e.g., ==), and comments was complex.
- Impact: Misclassified tokens broke parsing in ast.cpp, and comments were sometimes tokenized.
- Solution: Added context-aware tokenization for int main(), tracked previous characters for operators, and used removeAllcomments to skip comments:

```
if (str == "int" && next_str == "main" && next_ch == '(') {
```

```
        token[token1].name = "int_main";

        token[token1].type = "main_func";

        token[token1].id = MAIN_FUNC;

        token1++;

    }
```

These challenges—parsing complexity, memory management, and tokenization accuracy—were tackled with robust parsing, memory cleanup, and precise token handling, ensuring the tool's reliability. Future improvements can build on these solutions.


**7.Conclusion**

The AST Builder: C Edition project has successfully delivered a functional tool that processes C source code through a three-phase pipeline—tokenization, Abstract Syntax Tree (AST) construction, and visualization—enabling detailed analysis of C programs. This endeavor has provided a practical platform for understanding program structure, serving as a valuable resource for developers, researchers, and students in the realm of compiler design and program analysis.

The development of this project offered significant insights into both technical and practical aspects of software engineering:

- Compiler Design Principles: I deepened my understanding of lexical and syntactic analysis by implementing a custom tokenizer in tokenization.cpp and a recursive descent parser in ast.cpp. This hands-on experience clarified how source code is transformed into tokens and then into a structured AST, reinforcing concepts like grammar rules and tree traversal.
- C and C++ Proficiency: Working with C in ast.cpp and C++ in main.cpp and tokenization.cpp enhanced my skills in both languages. I learned to manage memory manually in C using malloc and free, ensuring no leaks in the AST, and utilized C++'s standard library for efficient file I/O and string handling, addressing cross-language compatibility challenges.
- Error Handling and Usability: Designing a user interface in main.cpp with colored output taught me the importance of clear feedback for usability. Implementing robust error handling—such as file I/O errors in tokenization.cpp and parsing errors in ast.cpp—highlighted the need for informative error messages to aid debugging and improve user experience.
- Modular Design and Integration: The project's modular structure, with separate components for tokenization, parsing, and orchestration, emphasized the benefits of separation of concerns. However, integrating these via system calls in main.cpp revealed challenges in ensuring binary availability and handling execution failures, teaching me the importance of robust integration strategies.
- Testing and Debugging: Testing with diverse C programs, including edge cases, and using tools like Valgrind to detect memory leaks in ast.cpp improved my debugging skills. This process underscored the value of comprehensive testing to ensure reliability and correctness.

The AST Builder: C Edition lays a strong foundation for future enhancements to expand its functionality and applicability:

- Support for Advanced C Features: Extend the parser in ast.cpp to handle pointers, structs, and preprocessor directives, enabling analysis of more complex C programs.
- Graphical User Interface (GUI): Replace the terminal-based interface with a GUI using frameworks like Qt, offering interactive AST visualization with features like node highlighting and zooming, enhancing usability for educational purposes.
- Semantic Analysis: Add a semantic analysis phase to check for type mismatches and undeclared variables, transforming the tool into a static analyzer with broader applications.
- Performance Optimization: Address the 10,000-token limit by implementing dynamic resizing of token arrays in tokenization.cpp and ast.cpp, improving scalability for large files.
- IDE Integration: Integrate the tool into IDEs like Visual Studio Code as a plugin, providing real-time AST visualization and analysis during development.
- Multi-Language Support: Extend the tool to support languages like C++ or Python by adapting the tokenizer and parser, potentially using a parser generator like ANTLR for efficiency.
- Enhanced Error Recovery: Improve error handling in ast.cpp to recover from syntax errors, allowing partial AST construction and better diagnostics for malformed code.

The AST Builder: C Edition project has been an enriching experience, deepening my knowledge of compiler design, C/C++ programming, and software engineering practices. The challenges faced and solutions implemented have provided a solid foundation for future work, and the proposed extensions offer exciting opportunities to enhance the tool's capabilities, making it a more powerful resource for code analysis and education in the field of software development.

## 8. Reference

[1] Aho, A. V., Lam, M. S., Sethi, R., & Ullman, J. D. *Compilers: Principles, Techniques, and Tools (2nd ed.)*. Addison-Wesley, 2006.

- URL:https://www.pearson.com/store/p/compilers-principles-techniques-and-tools/P100000252052
- Last accessed on: 19 Mar 2025

[2] Grune, D., & Jacobs, C. J. H. "Modern Parsing Techniques for Programming Languages." *Journal of Computer Science and Technology*, 2008, pp. 45-60.

- This academic paper, published in the *Journal of Computer Science and Technology*, discusses modern parsing techniques, including recursive descent parsing used in ast.cpp.

[3] Kernighan, B. W., & Ritchie, D. M. *The C Programming Language (2nd ed.)*. Prentice Hall, 1988.

- URL:https://www.pearson.com/store/p/c-programming-language/P100000252053
- Last accessed on: 19 Mar 2025