

# Assignment 2: Syntax Analysis with Bison

CS 464/5607: Compiler Design

Spring 2026

## 1 Introduction

In this assignment, you will implement the second phase of a compiler: the **Syntax Analyzer** (or Parser).

Your task is to use **Bison** to generate a parser that consumes the stream of tokens produced by the lexer and verifies if they follow the grammatical rules of our language. Furthermore, you will construct an **Abstract Syntax Tree (AST)** that represents the hierarchical structure of the source code.

## 2 Provided Files

The following files are provided in the skeleton code.

### 2.1 Lexer Source

- `lib/lexer.l`: This file contains the Lexical Analyzer definitions using Flex. During the build process, this file is compiled to generate the C code that handles token generation.

### 2.2 Source Files (Your Task)

- `src/parser.y`: The Bison grammar file. This is where you will define the grammar rules, precedence, and AST construction logic.
- `src/ast.c`: The implementation file for the Abstract Syntax Tree. You must implement the node creation and management functions here.
- `src/ast.h`: The header file defining the `ASTNode` structure and `NodeType` enums.

### 2.3 Driver Code (Do Not Modify)

- `src/driver.c`: The entry point that calls `yyparse()` and prints the resulting AST or error messages.

## 3 The Abstract Syntax Tree (AST)

An **Abstract Syntax Tree (AST)** is a tree representation of the abstract syntactic structure of source code. Unlike a Parse Tree, which contains every detail of the grammar (including parentheses and semicolons), an AST contains only the structural details essential for analysis and code generation.

### 3.1 Why do we need it?

The parser validates the code, but simply returning "Success" or "Failure" is insufficient for a compiler. We need a data structure that captures the meaning of the program so that subsequent phases (Semantic Analysis, Code Generation) can traverse it.

### 3.2 API Functions

You must implement the following functions in `src/ast.c`:

### 3.2.1 create\_node

```
1 ASTNode* create_node(NodeType type, char* data, int line, ASTNode* left,
ASTNode* right, ASTNode* extra);
```

Allocates memory for a new `ASTNode`, initializes its fields, and returns a pointer to it.

- `type`: The kind of node (e.g., `NODE_VAR_DECL`, `NODE_IF`).
- `data`: A string for storing values like identifiers ("x"), operators ("+"), or literals ("10").

### 3.2.2 append\_node

```
1 ASTNode* append_node(ASTNode* head, ASTNode* new_node);
```

Used for handling lists of statements or declarations. This function should attach `new_node` to the end of the `next` chain starting at `head`.

### 3.2.3 free\_ast

```
1 void free_ast(ASTNode* node);
```

Frees all memory associated with the AST nodes to prevent memory leaks.

### 3.2.4 Note on print\_ast

The function `void print_ast(ASTNode* node, int level)` has already been implemented for you. **DO NOT MODIFY THIS FUNCTION.** It formats the output exactly as required by the test cases.

## 4 The Parser (`parser.y`)

### 4.1 The %union

Bison parsers use a union to hold different types of semantic values. In your skeleton code, it is defined as:

```
1 %union {
2     char* sval;           // For lexemes (IDs, Literals)
3     struct ASTNode* node; // For AST Nodes
4 }
```

- `sval`: Used for tokens that carry raw text, such as `T_ID` (identifiers) or `T_INT_LIT` (numbers).
- `node`: Used for non-terminals (like `program`) that produce an AST node as a result.

#### Example Usage:

```
1 %token <sval> T_ID
2 %type <node> program
```

### 4.2 Operator Precedence and Associativity

Ambiguities in expressions (like `1 + 2 * 3`) must be resolved using precedence rules. In Bison, you define these using `%left`, `%right`, or `%nonassoc`.

- Operators declared **later** have **higher** precedence.
- `%left` indicates left-associativity (e.g., `a - b - c` becomes `(a - b) - c`).

#### Example:

```

1  /* Low Precedence */
2  %left T_PLUS T_MINUS
3  %left T_MULT T_DIV
4  /* High Precedence */

```

## 4.3 Grammar Rules

You must define the grammar rules for the language constructs. A Bison rule defines a non-terminal (LHS) in terms of terminals and other non-terminals (RHS), followed by a block of C code that executes when the rule is matched.

```

1 LHS : RHS_1 RHS_2 ... RHS_N { Action Code }

```

### 4.3.1 Example: Assignment Statement

To help you understand how to construct the AST, let's look at how you might implement an **assignment** rule (e.g., `x = 5;`).

```

1 assignment
2   : T_ID T_ASSIGN expression T_SEMI {
3     /* $1 corresponds to T_ID (the variable name, e.g., "x")
4      $3 corresponds to expression (the value, e.g., "5")
5      */
6
7     // 1. Create a node for the variable on the Left Hand Side
8     ASTNode* varNode = NEW_NODE(NODE_VAR_USE, $1, NULL, NULL, NULL);
9
10    // 2. Create the main Assignment node
11    //      Left Child = variable, Right Child = expression
12    $$ = NEW_NODE(NODE_ASSIGN, NULL, varNode, $3, NULL);
13  }
14 ;

```

#### Explanation of Symbols:

- **\$\$** (Dollar-Dollar): This represents the **result** of the current rule. In the example above, we assign the new `NODE_ASSIGN` pointer to **\$\$**, which passes it up to the parent rule (e.g., `statement`).
- **\$1**: This refers to the value of the **first** component on the RHS. Since `T_ID` is defined as `<sval>` in the `%union`, **\$1** holds the string name of the identifier (e.g., "count").
- **\$3**: This refers to the **third** component, which is `expression`.

### 4.3.2 Implementation Hints

- **Program Structure:** A program is a list of declarations (variables or functions).
- **Expressions:** Use the precedence rules defined above to handle math logic without writing complex recursive grammar rules manually.
- **Actions:** Inside the curly braces `{}`, use the `create_node` function to build the tree.
  - **\$\$** refers to the result of the LHS (the parent node).
  - **\$1**, **\$2**, etc., refer to the values of the symbols on the RHS.

## 4.4 Node Configuration Reference

When calling `create_node` inside your grammar actions, you must map the children pointers (`left`, `right`, `extra`) correctly to match the expected AST structure.

- **Variable Declaration:**

- **left**: The Type node (e.g., "int", "float").
- **right**: The Initialization expression (or NULL if not initialized).
- **extra**: The Array Size node if it is an array, otherwise NULL.

- **If Statement:**

- **left**: The Condition expression.
- **right**: The "Then" Block.
- **extra**: The "Else" Block (or NULL if no else exists).

- **While Loop:**

- **left**: The Condition expression.
- **right**: The Body Block.

- **For Loop:**

- **Syntax**: `for (init; condition; update;) body`
- **left**: The Initialization statement.
- **right**: The Condition expression.
- **extra**: The Body Block (including the update statement)

*Note:* Notice the semicolon at the end of the update statement. This simplifies the implementation by allowing the grammar to reuse the standard Assignment Statement rule directly, rather than requiring a separate rule for assignments without semicolons.

- **I/O Statements:**

- **data**: Set to "read" for `cin` or "print" for `cout`.
- **left**: The variable (for read) or expression (for print).

## 4.5 Panic Mode Error Recovery

To make your parser robust, you must implement **Panic Mode Recovery**. When a syntax error occurs, the parser should not crash; instead, it should discard tokens until it finds a "synchronizing token" (like a semicolon or a closing brace) and then resume parsing.

**How to implement:** Use the special `error` token provided by Bison.

```

1  statement
2    : ... (valid rules) ...
3    | error T_SEMI { yyerrok; }
4    ;

```

This tells Bison: "If you find an error, skip everything until you see a semicolon, then clear the error flag (`yyerrok`) and continue."

## 4.6 Error Reporting

The function `void yyerror(const char *s)` is already provided in the skeleton code. It is automatically called by Bison when a parse error occurs. It prints the line number and the current token to `stderr`. This has already been implemented in `parser.y` but you may modify it accordingly.

# 5 Building and Testing

## 5.1 Compilation

To build the project, open your terminal in the project root directory and run:

```

1  make

```

This will compile your parser source code and link it with the provided `lexer.o` object file. The executable will be placed in the `build/` directory.

## 5.2 Running the Test Suite

We have provided a Python script to automate testing. To run it, execute:

```
1 python run_tests.py
```

The script compares your parser's output against "Golden Output" files.

- **70%** of the grade is based on the visible tests provided to you.
- **30%** of the grade is based on **Hidden Test Cases** that check for edge cases and robustness.

## 5.3 Testing with Custom Inputs

You are encouraged to create your own test cases to debug specific issues.

1. Create a text file (e.g., `my_input.txt`) with some code.
2. Run your built parser and redirect the input:

**On Windows:**

```
1 build\parser.exe my_input.txt
```

**On macOS/Linux:**

```
1 ./build/parser my_input.txt
```

This will print the AST to the terminal, allowing you to manually verify if your rules are working as expected.

## 5.4 Platform Specifics (Windows)

If you are working on Windows natively (Command Prompt or PowerShell, not WSL), you must modify the build scripts to handle file extensions correctly.

### 1. Modify the Makefile

Windows executables require the `.exe` extension. Find the `TARGET` variable and update it:

```
1 # Change:
2 TARGET = $(BUILD_DIR)/parser
3
4 # To:
5 TARGET = $(BUILD_DIR)/parser.exe
```

### 2. Modify run\_tests.py

The Python script needs to know it should look for an executable with an extension.

```
1 # Change:
2 PARSER_EXEC = os.path.join("build", "parser")
3
4 # To:
5 PARSER_EXEC = os.path.join("build", "parser.exe")
```

## 6 Submission Requirements