# CFG for Minimal TypeScript

## CS F365 – Compiler Construction

Phase 1: Language Design and Lexical Analysis

## Grammar Notation

- *Uppercase* = Non-terminals
- `monospace` = Terminals (tokens)
- | = alternatives, $\varepsilon$ = empty production

## Start Symbol & Program Structure

$$Program \rightarrow StmtList$$

$$StmtList \rightarrow Stmt\ StmtList \mid \varepsilon$$

$$Stmt \rightarrow VarDecl$$
$$\mid AssignStmt$$
$$\mid IfStmt$$
$$\mid WhileStmt$$
$$\mid Block$$

$$Block \rightarrow \{\ StmtList\ \}$$

## Variable Declaration

$$VarDecl \rightarrow VarKeyword\ \texttt{id}\ \texttt{:}\ Type\ \texttt{=}\ Expr\ \texttt{;}$$
$$\mid VarKeyword\ \texttt{id}\ \texttt{:}\ Type\ \texttt{;}$$

$$VarKeyword \rightarrow \texttt{let} \mid \texttt{const}$$

$$Type \rightarrow \texttt{number} \mid \texttt{string} \mid \texttt{boolean}$$

## Assignment Statement

$$AssignStmt \rightarrow \texttt{id}\ \texttt{=}\ Expr\ \texttt{;}$$

## Conditional Statement (if-else)

$$IfStmt \rightarrow \texttt{if}\ \texttt{(}\ Expr\ \texttt{)}\ Block$$
$$\mid \texttt{if}\ \texttt{(}\ Expr\ \texttt{)}\ Block\ \texttt{else}\ Block$$
$$\mid \texttt{if}\ \texttt{(}\ Expr\ \texttt{)}\ Block\ \texttt{else}\ IfStmt$$

## Loop Statement (while)

$$WhileStmt \rightarrow \texttt{while (} Expr \texttt{)} Block$$

## Expressions (with Precedence)

Precedence increases from top to bottom (lowest → highest).

| | | |
|---|---|---|
| $Expr \rightarrow Expr$ `||` $AndExpr$ | | ← Logical OR (lowest |
|     \| $AndExpr$ | | |
| $AndExpr \rightarrow AndExpr$ `&&` $RelExpr$ | | ← Logical AND |
|     \| $RelExpr$ | | |
| $RelExpr \rightarrow RelExpr$ `relop` $AddExpr$ | | ← Relational (==, !=, |
|     \| $AddExpr$ | | |
| $AddExpr \rightarrow AddExpr$ `+` $MulExpr$ | | ← Addition |
|     \| $AddExpr$ `-` $MulExpr$ | | ← Subtraction |
|     \| $MulExpr$ | | |
| $MulExpr \rightarrow MulExpr$ `*` $UnaryExpr$ | | ← Multiplication |
|     \| $MulExpr$ `/` $UnaryExpr$ | | ← Division |
|     \| $MulExpr$ `%` $UnaryExpr$ | | ← Modulo |
|     \| $UnaryExpr$ | | |
| $UnaryExpr \rightarrow$ `!` $UnaryExpr$ | | ← Logical NOT |
|     \| `-` $UnaryExpr$ | | ← Unary minus |
|     \| $Primary$ | | |

$$Primary \rightarrow \texttt{id} \mid \texttt{number\_lit} \mid \texttt{string\_lit} \mid \texttt{true} \mid \texttt{false} \mid \texttt{(} Expr \texttt{)}$$

## Terminal Symbols Summary

| Category | Terminals |
|---|---|
| Keywords | `let`, `const`, `if`, `else`, `while`, `true`, `false` |
| Types | `number`, `string`, `boolean` |
| Operators | `+`, `-`, `*`, `/`, `%`, `&&`, `||`, `!` |
| Relational | `==`, `!=`, `<`, `>`, `<=`, `>=` |
| Delimiters | `(`, `)`, `{`, `}`, `:`, `=`, `;` |
| Literals | `number_lit`, `string_lit` |
| Identifier | `id` |

## Lex Priority Order

Lex resolves conflicts by two rules — **longest match first**, then **order of appearance**.
The recommended rule ordering is:

1. Multi-character operators    (`==`, `!=`, `<=`, `>=`, `&&`, `||`)
2. Keywords    (`let`, `const`, `if`, `else`, `while`, `true`, `false`)
3. Type keywords    (`number`, `string`, `boolean`)
4. Literals    (`number_lit`, `string_lit`)
5. Identifiers    (`id`)
6. Single-character operators    (`+`, `-`, `*`, `/`, `%`, `!`, `<`, `>`, `=`)
7. Delimiters    (`(`, `)`, `{`, `}`, `:`, `;`)
8. Whitespace & Comments    (skip)

## How to Run the Lexer

### Prerequisites

Ensure the following tools are installed on your system:

- `flex` — the lexer generator (implements Lex)
- `gcc` — the GNU C compiler

On Ubuntu/Debian, install them with:

```
sudo apt-get install flex gcc
```

### Step 1: Generate C Code from the Lexer

Run `flex` on the `.l` file to produce `lex.yy.c`:

```
flex typescript_lexer.l
```

### Step 2: Compile the Generated C Code

Compile `lex.yy.c` with `gcc` and link the flex library (`-lfl`):

```
gcc lex.yy.c -o ts_lexer -lfl
```

This produces the executable `ts_lexer`.

### Step 3: Run on Test Files

**Valid input program:**

```
./ts_lexer valid_input.ts
```

**Input program with lexical errors:**

```
./ts_lexer error_input.ts
```

**Alternatively, pipe input from stdin:**

```
./ts_lexer < valid_input.ts
```

**Expected Output Format**

The lexer prints a formatted token stream to stdout:

```
================================================================
         Minimal  TypeScript  Lexer   --   Token  Stream
================================================================
LINE    TOKEN TYPE                 LEXEME
----------------------------------------------------------------
4       KW_LET                     let
4       IDENTIFIER                 x
4       COLON                      :
4       TYPE_NUMBER                number
4       OP_ASSIGN                  =
4       NUM_LIT                    10
4       SEMICOLON                  ;
...
================================================================
  Total tokens  : 87
  Lexical errors: 0
================================================================
```

For the error input, unrecognized characters are reported inline:

```
7       *** LEXICAL ERROR ***      Unexpected character: '@'
10      *** LEXICAL ERROR ***      Unexpected character: '#'
...
  Lexical errors: 4
```

The lexer returns exit code 0 on success and 1 if any lexical errors were found, making it suitable for use in automated build pipelines.