

# 编译原理实验报告

|      |                                      |
|------|--------------------------------------|
| 实验1  | 词法分析器 (Lexical Analyzer Programming) |
| 实验2  | 语法分析器 (Syntax Parser)                |
| 姓名   | 熊宇祺                                  |
| 学号   | 09023112                             |
| 报告日期 | 2025年12月28日                          |

## 实验一：词法分析器

### a) 实验动机 / 目的

- 理解词法分析原理：** 深入理解编译器前端的工作机制，特别是如何将源代码字符流转换为抽象的词法单元 (Token) 序列。
- 掌握有限自动机 (FA) 应用：** 学习如何将正则表达式 (REs) 转化为确定有限自动机 (DFA)，并使用编程语言 (C语言) 实现基于 DFA 的状态跳转逻辑。
- 提升编程实践能力：** 练习文件 I/O 操作、字符串处理以及缓冲区管理，为后续的语法分析实验打下基础。

### b) 内容描述

本实验旨在实现一个功能完备的词法分析器 (Scanner)。

**输入：** 包含语言源代码的文本文件。

**输出：** 识别出的 Token 序列 (Tag + Attribute)，并打印至控制台及文件。

#### 词法单元定义 (Token Definitions):

本实验针对 TINY 语言的一个子集，需要识别以下具体的词法单元类型：

#### 1. Reserved Words (保留字)

- 集合：** if, then, else, end, repeat, until, read, write
- 特性：** 大小写敏感 (Case-sensitive)。即 if 是关键字，而 IF 或 If 将被识别为标识符 (ID)。

#### 2. Identifiers & Numbers (自定义符)

- ID (标识符)：** 用于变量命名。以字母开头，后跟任意数量的字母或数字。
- NUM (整数)：** 无符号整数序列。由连续的数字字符组成。

#### 3. Special Symbols & Operators (特殊符号与运算符)

- ASSIGN (赋值)：** :=
- RELOP (关系运算符)：** < (小于), = (等于) 扩充支持：根据实验描述，同时支持 > (大于), <= (小于等于), >= (大于等于), <> (不等于)。
- ARITHOP (算术运算符)：** +, -, \*, /
- PUNCTUATION (标点)：** ( (左括号), ) (右括号), ; (分号)

#### 4. Filtered Elements (过滤元素)

- COMMENT (注释):** 以 { 开始, 以 } 结束的文本块。注释内容不生成 Token。
- WHITESPACE (空白):** 空格、制表符 \t、换行符 \n。用于分隔 Token, 通过 skip 操作处理。

## c) 设计思路 / 方法

### 1. 总体设计思路

本实验遵循 "Programming based on FA" 方法论, 具体步骤如下:

- 正则定义 (RE Definition):** 形式化描述每一类 Token 的字符构成规则。
- 自动机建模 (Modeling):** 将 RE 转化为 NFA (非确定有限自动机)。将 NFA 合并并转换为 DFA (确定有限自动机)。对 DFA 进行最小化处理, 得到最终的状态转换图。
- 代码映射 (Implementation):** 使用 C 语言的 switch-case 结构直接模拟最小化 DFA 的状态跳转逻辑。
- 混合策略 (Hybrid Strategy):** 针对关键字识别, 采用 "先识别为 ID, 再查表 (Lookup)" 的策略, 以简化 DFA 的状态数量。

### 2. 正则表达式定义

#### 2.1 关键字正则表达式

```
IF:      (if|IF)
THEN:    (then|THEN)
ELSE:    (else|ELSE)
END:     (end|END)
REPEAT:  (repeat|REPEAT)
UNTIL:   (until|UNTIL)
READ:    (read|READ)
WRITE:   (write|WRITE)
```

#### 2.2 标识符正则表达式

```
ID: [a-zA-Z_][a-zA-Z0-9_]*
```

#### 2.3 数字正则表达式

```
NUM: [0-9]+
```

#### 2.4 运算符正则表达式

```
ASSIGN: :=
RELOP:  (<>|>=|<=|>|<)
PLUS:   \+
MINUS:  -
TIMES:  \*
OVER:   /
```

## 2.5 分隔符正则表达式

```
LPAREN: \(
RPAREN: \)
SEMI: ;
DELIMITER: [\t\n\r ]+
```

## 3. 有限自动机构造过程

### 1. Thompson 构造法 (RE → NFA)

Thompson 构造法用于将正则表达式 (RE) 系统性地转换为不确定有限自动机 (NFA)。

- 核心思想：将复杂的正则表达式拆解为若干基本构件，每个基本构件对应一个最小的 NFA 片段，再通过  $\epsilon$ -transition 将这些片段组合起来。
- 构造原则如下：
  - 每一个基本符号 (如 `digit`、`letter`) 都对应一个简单的 NFA；
  - 正则表达式中的连接、选择和闭包操作，分别通过 NFA 片段的顺序连接、分支合并和自循环来实现；
  - 组合过程中大量使用  $\epsilon$ -transition，使得状态之间无需消耗输入字符即可转移。

例如，对于整数常量的正则表达式：

```
NUM → digit+
```

### 2. 子集构造法 (NFA → DFA)

由于 NFA 在同一输入符号下可能存在多条转移路径，实际词法分析器实现中通常需要将 NFA 转换为确定有限自动机 (DFA)。子集构造法 (Subset Construction) 正是完成这一转换的标准方法。

基本思想：**DFA 中的每个状态，都对应 NFA 中的一组状态（一个状态集合）。**

基本步骤如下：

#### (1) 引入 $\epsilon$ -transition

- 将各个 Token 对应的 NFA 的起始状态，通过  $\epsilon$ -transition 连接到一个统一的 `start` 节点；
- 该 `start` 节点作为整个词法分析器的全局入口状态；
- 通过这种方式，可以在同一个自动机中同时识别多种 Token。

#### (2) 计算 $\epsilon$ -closure

- 对 `start` 节点计算其  $\epsilon$ -closure；
- $\epsilon$ -closure 表示：从某一状态出发，仅通过  $\epsilon$ -transition 能够到达的所有状态集合；
- 该  $\epsilon$ -closure 集合即作为 **DFA 的初始状态**。

#### (3) 状态扩展与转移计算

- 对 DFA 中的每一个状态集合 `T`，以及每一个输入符号 `input`，执行以下操作：

- 计算  $\text{move}(T, \text{input})$ ，得到在输入符号 `input` 下，从集合 `T` 中任一状态出发能够到达的 NFA 状态集合；
- 对  $\text{move}(T, \text{input})$  的结果再计算其  $\epsilon$ -closure；
- 将得到的  $\epsilon$ -closure 集合作为一个新的 DFA 状态。

#### (4) 构造完整的状态转移表

- 将新生成的 DFA 状态加入状态集合；
- 对尚未处理过的 DFA 状态重复步骤 (3)；
- 直到不再产生新的 DFA 状态为止；
- 最终即可得到一个完整、确定的 DFA 及其状态转移表。

### 3. 冲突解决策略：最长匹配原则

在词法分析过程中，不同的 Token 可能具有相同的前缀，从而导致识别冲突。为保证词法分析结果的正确性，本实验采用 **最长匹配原则 (Maximal Munch)** 进行冲突解决。

- 基本思想：在所有能够匹配当前输入串的 Token 中，**优先选择匹配字符数最多的那个**。

以关系运算符的识别过程为例：

- 当扫描到字符 `<` 时，DFA 不会立即接受该 Token；
- 扫描器会继续向前预读一个字符以判断是否能够形成更长的匹配：
  - 若下一个字符为 `=`，则识别为关系运算符 `<=`；
  - 否则，将多读入的字符回退，并接受 `<` 作为最终的 Token。

通过遵循最长匹配原则，可以有效解决 Token 之间由于公共前缀引发的冲突问题，从而确保词法分析器在复杂输入情况下仍能产生唯一且正确的 Token 序列。

### 4. 实现方法选择

- **缓冲区管理**：使用双缓冲区或单行缓冲区策略，利用 `ungetc()` 处理超前扫描 (Lookahead) 带来的字符回退需求。
- **DFA 模拟方式**：采用 **直接编码法 (Direct Coded Approach)**。相比于查表驱动法 (Table-Driven)，直接利用控制流 (Switch-Case) 在处理复杂的副作用 (如注释内的换行计数、错误恢复) 时更加灵活。

## d) 假设条件

1. **字符集**：源代码仅包含 ASCII 字符。
2. **注释格式**：假设注释以 `{` 开始，以 `}` 结束 (根据代码实现)，且不支持嵌套注释。
3. **标识符限制**：标识符由字母开头，后跟字母或数字，长度不超过 99 个字符。
4. **大小写敏感**：关键字 (如 `if`) 是大小写敏感的 (即 `IF` 会被识别为标识符而非关键字)。
5. **数字格式**：仅处理无符号整数 (代码逻辑主要针对整数，遇到 `.` 会停止或报错，除非扩充浮点逻辑)。

## e) 有限自动机设计与演化

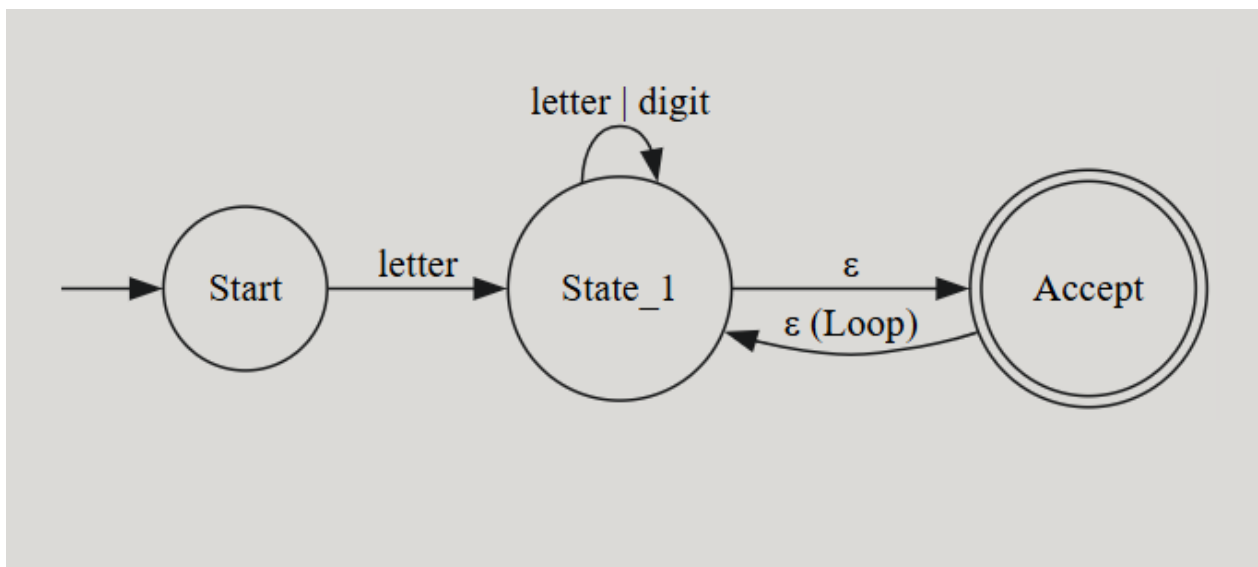
本实验严格遵循编译原理中的词法分析器构造流程：

- 正则表达式 (RE) → NFA → 合并 NFA → 最小化 DFA → 代码实现。以下是每个阶段的详细推导过程。

1. 定义正则表达式：前文已经定义过，不再赘述

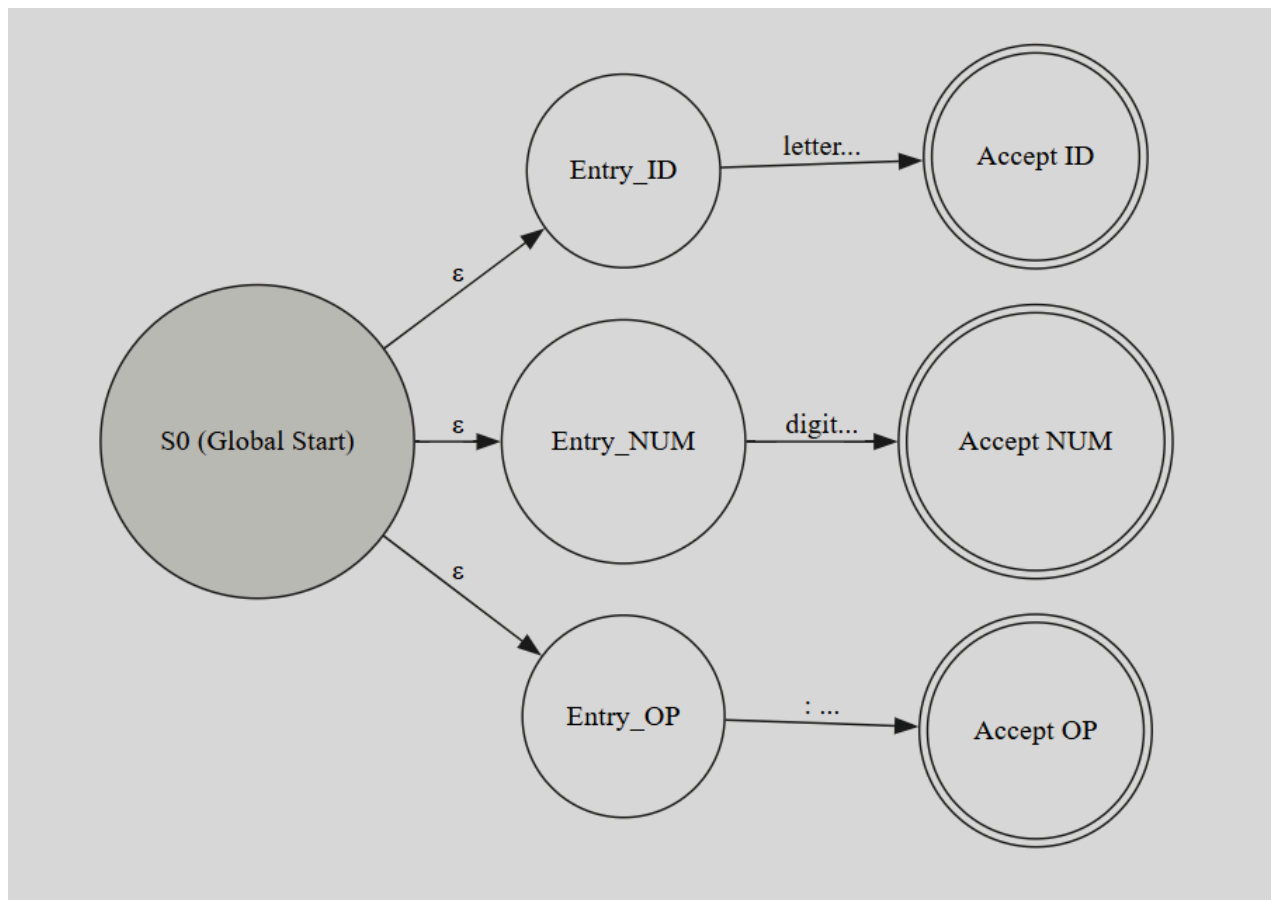
2. RE 转换 NFA：

- 利用 **Thompson 构造法**，将上述 RE 分解为基本的非确定有限自动机 (NFA) 片段。这里以最典型的 **标识符 (ID)** 为例进行展示。



3. 合并 NFA：

- 为了构建统一的词法分析器，我们引入一个新的全局起始状态 **s0 (Start)**。利用  $\epsilon$ -转移 (Epsilon Transitions) 将 **S0** 连接到所有子 NFA 的起始节点。这就构成了一个能非确定地识别所有 Token 的大型 NFA。



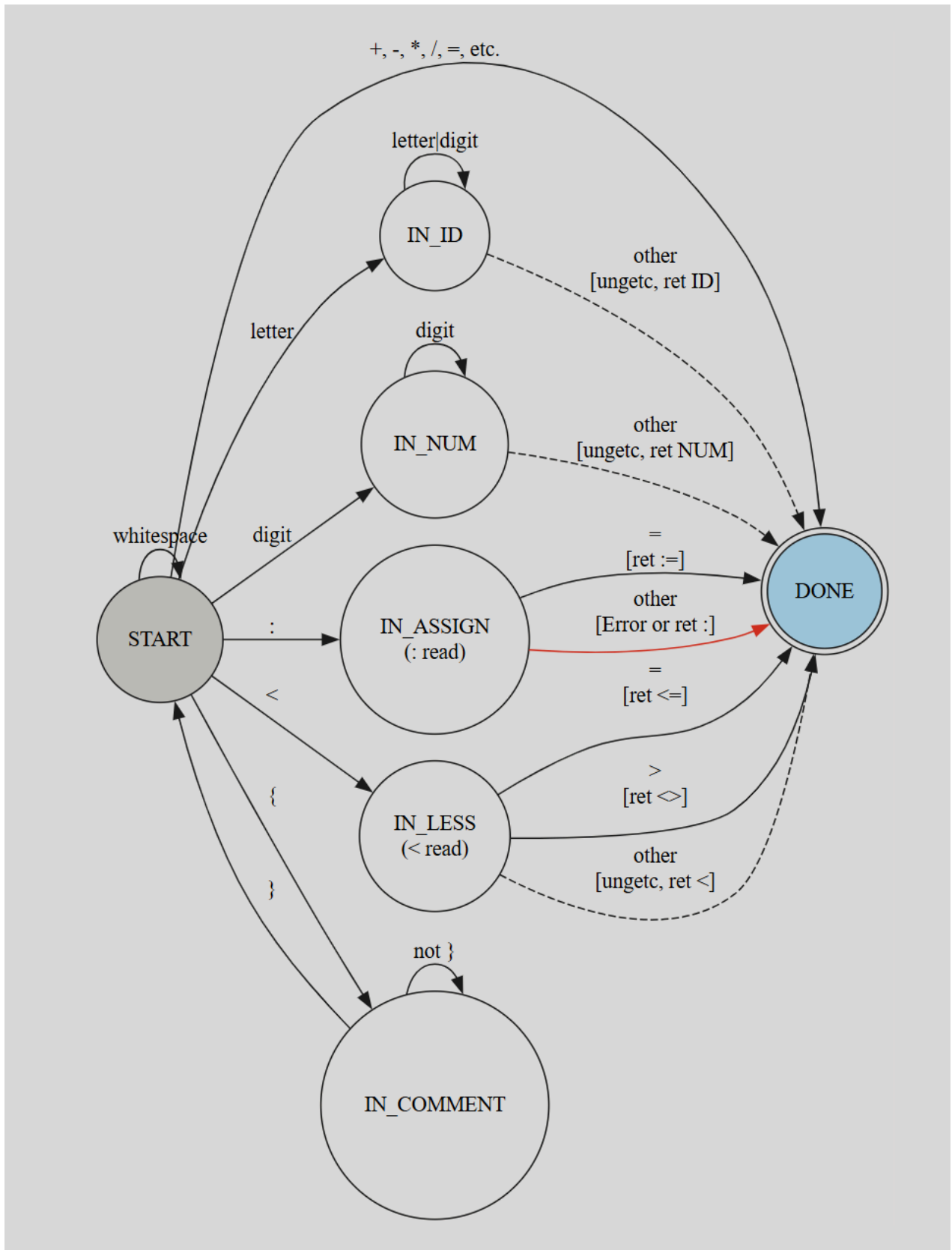
#### 4. NFA 转 DFA 并最小化:

通过 **子集构造法 (Subset Construction)** 消除非确定性，并处理前缀冲突（如 `<` 和 `<=`）。随后通过 **Hopcroft 算法** 或等价类划分进行最小化，得到 **DFA'**。

这个 DFA' 是我们编程实现的直接蓝图（State Diagram）。

- **状态说明:**

- **START:** 初始状态，等待输入。
- **IN\_ID / IN\_NUM:** 正在读取标识符或数字。
- **IN\_ASSIGN:** 已读入 `:`，等待 `=`。
- **IN\_LESS:** 已读入 `<`，根据下一字符决定是 `<`、`<=` 还是 `<>`。
- **DONE:** 接受状态（代码中通过 `return` 实现）。



## 5. 基于 DFA' 的编程实现

在 C 语言实现中，我们将上述 DFA' 的状态节点映射为 StateType 枚举，将边映射为 switch-case 分支结构。

DFA 到 代码的映射表：

| DFA 状态节点   | C 语言状态              | 行为逻辑   |
|------------|---------------------|--|
| START      | case START:         | 读取字符 c，根据字符类型切换 state。                         |
| IN_ID      | case IN_ID:         | 若 c 是字母/数字，保持状态；否则 ungetc(c)，state=DONE，返回 ID。 |
| IN_ASSIGN  | case IN_ASSIGN:     | 若 c=='='，返回 ASSIGN；否则报错。                       |
| IN_LESS    | case IN_LESS:       | 若 c=='=' 返 LE；若 c=='>' 返 NE；否则 ungetc 并返 LT。   |
| IN_COMMENT | case<br>IN_COMMENT: | 循环读取直到 }，然后重置为 START（不生成 Token）。               |
| DONE       | while(state!=DONE)  | 循环终止条件。在此处进行关键字查表（Lookup）并返回最终 Token。          |

这种基于自动机的编程方式（Table-Driven 或 Direct-Coded）保证了词法分析器的高效性、可扩展性以及与理论模型的一致性。

## f) 重要数据结构说明

### 1. Token 类型枚举 (TokenType)

用于区分词法单元类别。我们在设计中将所有可能的输出结果映射为整数枚举值，便于语法分析器消费。

```
typedef enum {
    /* Bookkeeping Tokens (辅助标记) */
    ENDFILE, ERROR,

    /* Reserved words (保留字) */
    IF, THEN, ELSE, END, REPEAT, UNTIL, READ, WRITE,

    /* Multicharacter Tokens (多字符标记) */
    ID,      // 标识符
    NUM,     // 整数

    /* Special Symbols (特殊符号) */
    ASSIGN, // :=
    EQ,     // =
    LT,     // <
    LE,     // <=
    GT,     // >
    GE,     // >=
    NE,     // <>
    PLUS, MINUS, TIMES, OVER, LPAREN, RPAREN, SEMI
} TokenType;
```

### 2. DFA 状态枚举 (StateType)

这是实现“基于有限自动机编程”的核心。枚举值与我们在 **步骤 (d) 最小化 DFA** 中设计的状态节点一一对应。



```
typedef enum {
    START,           // 初始态: 准备接收新字符
    IN_ASSIGN,        // 赋值态: 已读入 ':', 等待 '='
    IN_COMMENT,       // 注释态: 已读入 '{', 忽略后续字符直到 '}'
    IN_NUM,           // 数字态: 正在读取数字序列
    IN_ID,            // 标识符态: 正在读取字母数字序列
    IN_LESS,          // 小于态: 已读入 '<', 处理 <=, <>, <
    IN_GREATER,       // 大于态: 已读入 '>', 处理 >=, >
    DONE              // 完成态: Token 识别结束, 准备返回
} StateType;
```

### 3. 保留字查找表 (ReservedWord Table)

为了简化 DFA 结构（避免为每个关键字建立独立的状态路径），我们采用 **"混合策略"**：先将关键字识别为普通标识符 (ID)，再通过查表法修正类型。

```
#define MAXRESERVED 8

typedef struct {
    char* str;           // 关键字字符串, 如 "if"
    TokenType tok;       // 对应的 Token 类型, 如 IF
} ReservedWord;

static ReservedWord reservedWords[MAXRESERVED] = {
    {"if", IF}, {"then", THEN}, {"else", ELSE}, {"end", END},
    {"repeat", REPEAT}, {"until", UNTIL}, {"read", READ}, {"write", WRITE}
};
```

### 4. 扫描缓冲区与全局变量

- char tokenString[MAXTOKENLEN+1]: 用于存储当前正在构建的 Token 的字符串值 (Lexeme)。
- int lineno: 记录当前源文件的行号，用于错误定位。
- FILE \* source: 指向源代码文件的指针。

## g) 核心算法说明

### 1. 主扫描算法

这是词法分析器的驱动引擎。它直接模拟 DFA 的状态转移逻辑。

**算法流程：**

1. **初始化：** 设置当前状态 state = START，清空 tokenString。
2. **DFA 循环：** 当 state != DONE 时重复执行：
  - **读取：** 获取下一个字符 c。
  - **保存判断：** 默认保存字符 c 到缓冲区，但对于空白符、注释内容或超前扫描失败回退的字符，标记 save = FALSE。
  - **状态转移 (Switch-Case)：**

- 若在 START: 根据 c 的类别跳转到 IN\_ID, IN\_NUM 等状态。
  - 若在 IN\_ID / IN\_NUM: 持续读取直到遇到非匹配字符, 执行 **回退 (Unget)**, 跳转至 DONE。
  - 若在 IN\_LESS (处理 <):
    - 若 c == '=' → 识别为 <=, state = DONE。
    - 若 c == '>' → 识别为 <>, state = DONE。
    - 否则 → 识别为 <, 执行 **回退**, state = DONE。
  - **缓冲**: 若 save == TRUE, 将 c 追加到 tokenString。
3. **关键字识别**: 循环结束后, 如果 Token 类型是 ID, 调用 reservedLookup() 检查是否为保留字。
4. **返回**: 返回最终识别的 TokenType。

## 2. 超前扫描与回退机制

为了实现**最大匹配原则 (Maximal Munch)**, 算法经常需要“多读一个字符”来判断当前 Token 是否结束。

- **场景**: 在识别数字 123 时, 必须读到第 4 个字符 (如空格或 +) 才能确定数字结束。
- **实现**:

```
case IN_NUM:
    if (!isdigit(c)) {
        ungetNextChar(); // 核心: 将多读的字符退回输入流
        save = FALSE;    // 该字符不属于当前数字
        state = DONE;    // 数字识别完成
        currentToken = NUM;
    }
    break;
```

ungetNextChar() 保证了下一个 Token 的识别能正确地该字符开始, 而不是将其丢失。

## 3. 保留字查找算法

在 DFA 将字符串识别为 ID 后, 执行此线性查找算法。

**算法逻辑**:

1. 接收输入字符串 s。
2. 遍历 reservedWords 数组。
3. 使用 strcmp 比较 s 与表中的关键字。
4. **匹配成功**: 返回表中对应的 TokenType (如 IF)。
5. **匹配失败**: 保持原类型, 返回 ID。
6. **注**: 由于 TINY 语言关键字很少 (8个), 线性查找效率足够高, 无需使用哈希表。

## 4. 注释过滤算法

注释的处理被内嵌在 DFA 的 IN\_COMMENT 状态中, 实现了在词法层面的自动过滤。

### 算法逻辑:

- 当在 START 遇到 { 时, 进入 IN\_COMMENT。
- 在 IN\_COMMENT 中, 持续读取字符但不保存 (save = FALSE)。
- 一旦遇到 }, 状态重置为 START 而不是 DONE。
- 效果: getToken 函数会继续循环, 就像从未见过这些注释字符一样, 直到找到一个有效的 Token 才返回。这确保了语法分析器永远不会看到注释。

## h) 实验核心代码

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <ctype.h>

// =====
// 1. 定义数据结构
// =====

// Token 类型定义
typedef enum {
    // 关键字
    IF, THEN, ELSE, END, REPEAT, UNTIL, READ, WRITE,
    // 类别
    ID, NUM,
    // 运算符
    PLUS, MINUS, TIMES, OVER,    // + - * /
    ASSIGN,                      // :=
    EQ, LT, GT, LTE, GTE, NEQ,   // = < > <= >= <>
    // 界符
    LPAREN, RPAREN, SEMI,        // ( ) ;
    // 结束与错误
    ENDFILE, ERROR
} TokenType;

// DFA 状态定义
typedef enum {
    START,
    IN_ID,
    IN_NUM,
    IN_ASSIGN,    // 读到了 :
    IN_LESS,      // 读到了 <
    IN_GREATER,   // 读到了 >
    IN_COMMENT,   // 读到了 {
    DONE          // 完成一个Token识别
} StateType;

// 全局变量
FILE *source;    // 输入源文件指针
FILE *outputFile; // 输出结果文件指针
```

```

char tokenString[100];
int lineNo = 1;

// 关键字查找表
struct {
    char *str;
    TokenType tok;
} keywords[] = {
    {"if", IF}, {"then", THEN}, {"else", ELSE}, {"end", END},
    {"repeat", REPEAT}, {"until", UNTIL}, {"read", READ}, {"write", WRITE},
    {NULL, 0}
};

// =====
// 2. 辅助函数
// =====

// 查表判断是否为关键字
TokenType lookup(char *s) {
    for (int i = 0; keywords[i].str != NULL; i++) {
        if (strcmp(s, keywords[i].str) == 0)
            return keywords[i].tok;
    }
    return ID;
}

char getNextChar() {
    return fgetc(source);
}

void ungetNextChar(char c) {
    ungetc(c, source);
}

// 使用 sprintf 先格式化，然后同时写入屏幕和文件
void printToken(TokenType token, const char *tokenString) {
    char buffer[256]; // 临时缓冲区，用于存放一行输出信息

    switch (token) {
        case IF: case THEN: case ELSE: case END:
        case REPEAT: case UNTIL: case READ: case WRITE:
            sprintf(buffer, "KEYWORD: %s\n", tokenString); break;
        case ASSIGN: sprintf(buffer, "ASSIGN: :=\n"); break;
        case LT: sprintf(buffer, "RELOP: <\n"); break;
        case GT: sprintf(buffer, "RELOP: >\n"); break;
        case EQ: sprintf(buffer, "RELOP: =\n"); break;
        case NEQ: sprintf(buffer, "RELOP: <>\n"); break;
        case LTE: sprintf(buffer, "RELOP: <=\n"); break;
        case GTE: sprintf(buffer, "RELOP: >=\n"); break;
        case LPAREN: sprintf(buffer, "SEMI: (\n"); break;
        case RPAREN: sprintf(buffer, "SEMI: )\n"); break;
        case SEMI: sprintf(buffer, "SEMI: ;\n"); break;
    }
}

```

```

        case PLUS: sprintf(buffer, "OP: +\n"); break;
        case MINUS: sprintf(buffer, "OP: -\n"); break;
        case TIMES: sprintf(buffer, "OP: *\n"); break;
        case OVER: sprintf(buffer, "OP: /\n"); break;
        case NUM: sprintf(buffer, "NUM: %s\n", tokenString); break;
        case ID: sprintf(buffer, "ID: %s\n", tokenString); break;
        case ENDFILE: sprintf(buffer, "EOF\n"); break;
        case ERROR: sprintf(buffer, "ERROR: Unexpected character '%s' at line %d\n",
tokenString, lineNo); break;
        default: sprintf(buffer, "UNKNOWN TOKEN\n"); break;
    }

    // 1. 输出到屏幕
    printf("%s", buffer);

    // 2. 输出到文件 (如果文件打开成功)
    if (outputFile != NULL) {
        fprintf(outputFile, "%s", buffer);
    }
}

// =====
// 3. 核心算法: DFA 驱动
// =====
TokenType getToken() {
    int tokenStringIndex = 0;
    TokenType currentToken;
    StateType state = START;
    int save;

    while (state != DONE) {
        char c = getNextChar();
        save = 1;

        switch (state) {
            case START:
                if (isdigit(c)) state = IN_NUM;
                else if (isalpha(c)) state = IN_ID;
                else if (c == ':') state = IN_ASSIGN;
                else if (c == '<') state = IN_LESS;
                else if (c == '>') state = IN_GREATER;
                else if (c == ' ' || c == '\t' || c == '\r') save = 0;
                else if (c == '\n') { save = 0; lineNo++; }
                else if (c == '{') { save = 0; state = IN_COMMENT; }
                else {
                    state = DONE;
                    switch (c) {
                        case EOF: save = 0; currentToken = ENDFILE; break;
                        case '=': currentToken = EQ; break;
                        case '+': currentToken = PLUS; break;
                        case '-': currentToken = MINUS; break;
                        case '*': currentToken = TIMES; break;
                        case '/': currentToken = OVER; break;
                    }
                }
            }
        }
    }
}

```

```

        case '(': currentToken = LPAREN; break;
        case ')': currentToken = RPAREN; break;
        case ';': currentToken = SEMI; break;
        default: currentToken = ERROR; break;
    }
}
break;

case IN_COMMENT:
    save = 0;
    if (c == '}') state = START;
    else if (c == '\n') lineNo++;
    else if (c == EOF) { state = DONE; currentToken = ENDFILE; }
    break;

case IN_ASSIGN:
    state = DONE;
    if (c == '=') currentToken = ASSIGN;
    else { ungetNextChar(c); save = 0; currentToken = ERROR; }
    break;

case IN_LESS:
    state = DONE;
    if (c == '=') currentToken = LTE;
    else if (c == '>') currentToken = NEQ;
    else { ungetNextChar(c); save = 0; currentToken = LT; }
    break;

case IN_GREATER:
    state = DONE;
    if (c == '=') currentToken = GTE;
    else { ungetNextChar(c); save = 0; currentToken = GT; }
    break;

case IN_NUM:
    if (!isdigit(c) && c != '.') {
        ungetNextChar(c); save = 0; state = DONE; currentToken = NUM;
    }
    break;

case IN_ID:
    if (!isalnum(c)) {
        ungetNextChar(c); save = 0; state = DONE; currentToken = ID;
    }
    break;

case DONE: default:
    state = DONE; currentToken = ERROR; break;
}

if ((save) && (tokenStringIndex < 99)) tokenString[tokenStringIndex++] = c;
if (state == DONE) {
    tokenString[tokenStringIndex] = '\0';
}

```

```

        if (currentToken == ID) currentToken = lookup(tokenString);
    }
}
return currentToken;
}

// =====
// 4. 主程序
// =====
int main() {
    // 1. 生成测试用的输入文件
    FILE *fp = fopen("test_code.txt", "w");
    if (fp) {
        fprintf(fp, "read x;\n");
        fprintf(fp, "if 0 < x then\n");
        fprintf(fp, "    fact := 1;\n");
        fprintf(fp, "    { Comment Ignored }\n");
        fprintf(fp, "    repeat fact := fact * x; until x = 0\n");
        fprintf(fp, "end");
        fclose(fp);
    }

    // 2. 打开输入文件
    source = fopen("test_code.txt", "r");
    if (source == NULL) {
        printf("Error: Could not open source file.\n");
        return 1;
    }

    // 3. 【修改】打开输出文件
    outputFile = fopen("output.txt", "w");
    if (outputFile == NULL) {
        printf("Error: Could not create output file.\n");
        // 即使输出文件创建失败，我们仍然继续运行，只是只输出到屏幕
    }

    printf("Analysis started. Results will be saved to 'output.txt'.\n");
    printf("=====\n");

    TokenType token;
    while ((token = getToken()) != ENDFILE) {
        printToken(token, tokenString);
    }

    // 4. 清理资源
    fclose(source);
    if (outputFile != NULL) {
        fclose(outputFile);
        printf("=====\n");
        printf("Analysis finished. Check 'output.txt' for results.\n");
    }

    return 0;
}

```

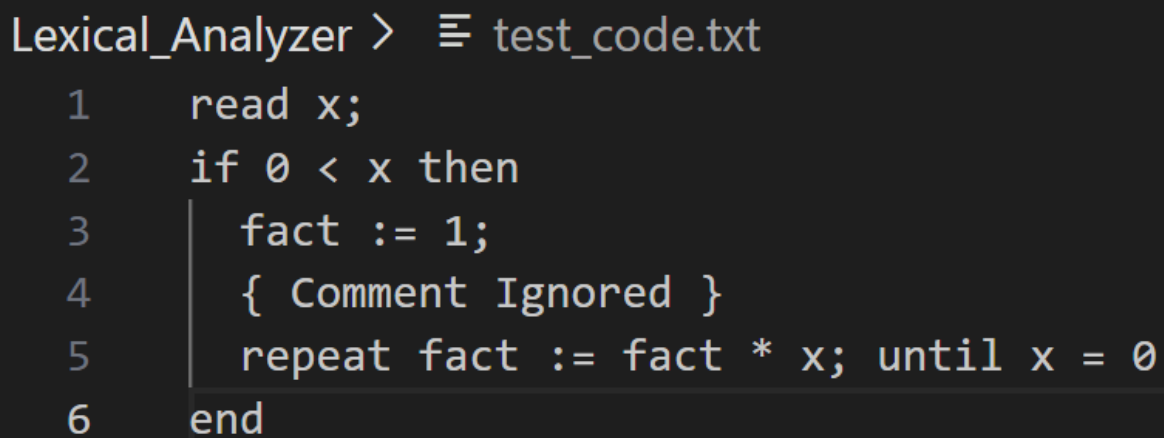
```
}
```

## i) 运行测试用例

输入文件 (test\_code.txt):

```
read x;  
if 0 < x then  
  fact := 1;  
  { Comment Ignored }  
  repeat fact := fact * x; until x = 0  
end
```

截图如下:



The screenshot shows a terminal window titled "Lexical\_Analyzer > test\_code.txt". It displays the source code from the previous block, with line numbers 1 through 6 on the left. The code is: 1 read x; 2 if 0 < x then 3 fact := 1; 4 { Comment Ignored } 5 repeat fact := fact \* x; until x = 0 6 end. The text is displayed in a monospaced font on a dark background.

```
Lexical_Analyzer > test_code.txt  
1 read x;  
2 if 0 < x then  
3   fact := 1;  
4   { Comment Ignored }  
5   repeat fact := fact * x; until x = 0  
6   end
```

运行结果 (截取自 output.txt):



Lexical\_Analyzer > ≡ output.txt

```
1  KEYWORD: read
2  ID: x
3  SEMI: ;
4  KEYWORD: if
5  NUM: 0
6  RELOP: <
7  ID: x
8  KEYWORD: then
9  ID: fact
10 ASSIGN: :=
11 NUM: 1
12 SEMI: ;
13 KEYWORD: repeat
14 ID: fact
15 ASSIGN: :=|
16 ID: fact
17 OP: *
18 ID: x
19 SEMI: ;
20 KEYWORD: until
21 ID: x
22 RELOP: =
23 NUM: 0
24 KEYWORD: end
25
```

(注: 程序成功跳过了{ Comment Ignored } 这一行注释)

| 行号 | Token 类型 (Tag) | 属性值 (Attribute/String) | 说明 (Description) |
|----|----------------|------------------------|------------------|
| 1  | KEYWORD        | read                   | 关键字 read         |
| 2  | ID             | x                      | 标识符 x            |
| 3  | SEMI           | ;                      | 分号               |
| 4  | KEYWORD        | if                     | 关键字 if           |
| 5  | NUM            | 0                      | 数字常量             |
| 6  | RELOP          | <                      | 小于号              |
| 7  | ID             | x                      | 标识符 x            |
| 8  | KEYWORD        | then                   | 关键字 then         |
| 9  | ID             | fact                   | 标识符 fact         |
| 10 | ASSIGN         | :=                     | 赋值符号             |
| 11 | NUM            | 1                      | 数字常量             |
| 12 | SEMI           | ;                      | 分号               |
| 13 | KEYWORD        | repeat                 | 关键字 repeat       |

| 行号 | Token 类型 (Tag) | 属性值 (Attribute/String) | 说明 (Description) |
|----|----------------|------------------------|------------------|
| 14 | ID             | fact                   | 标识符 fact         |
| 15 | ASSIGN         | :=                     | 赋值符号             |
| 16 | ID             | fact                   | 标识符 fact         |
| 17 | OP             | *                      | 乘法运算符            |
| 18 | ID             | x                      | 标识符 x            |
| 19 | SEMI           | ;                      | 分号               |
| 20 | KEYWORD        | until                  | 关键字 until        |
| 21 | ID             | x                      | 标识符 x            |
| 22 | RELOP          | =                      | 等于号              |
| 23 | NUM            | 0                      | 数字常量             |
| 24 | KEYWORD        | end                    | 关键字 end          |
| 25 | EOF            | -                      | 文件结束符            |

## j) 遇到的问题及解决方案

1. 问题：注释处理逻辑混淆
  - 描述：最初将注释视为一种特殊的 Token，导致输出中包含注释内容。
  - 解决：修改 DFA，在 IN\_COMMENT 状态下不保存字符（save = 0），并且在遇到结束符 } 后直接跳转回 START 状态，而不是 DONE 状态。
2. 问题：区分关键字与标识符
  - 描述：所有的关键字（如 if）本质上符合标识符的正则规则，很容易被误判为 ID
  - 解决：采用“先整体后局部”的策略。DFA 将所有字母开头的串都先识别为 ID，识别结束后，统一查表（Lookup Table）。如果在表中找到，则修正为对应的关键字类型，否则保留为 ID。
3. 问题：多读字符的处理
  - 描述：识别 NUM 或 ID 时，必须读到非数字/非字母才能停止，这导致文件指针多前进了一位。
  - 解决：使用标准库函数 ungetc() 实现回退机制，将多读的那个字符放回输入流，供下一次 getToken 使用。

## k) 心得体会

通过本次实验，我深刻体会到了有限自动机在计算机科学中的基础作用。虽然手动编写大量的 switch-case 看起来有些繁琐，但它清晰地展示了状态流转的过程，让我明白了编译器是如何“理解”源代码的第一步的。

这种基于 DFA 的设计具有良好的扩展性。如果将来需要增加新的运算符或语法规则（如支持浮点数或字符串字面量），只需在 switch 结构中增加相应的状态分支即可。此外，同时输出到屏幕和文件的功能也让我复习了 C 语言的文件操作，是一次非常有价值的实践。

# 实验二：语法分析器

## a) 实验动机 / 目的

- 理解语法分析原理**：深入理解编译器前端的核心组件——语法分析器的工作机制，掌握如何将线性的 Token 流转换为具有层级结构的语法树 (Syntax Tree)。
- 掌握 LL(1) 方法**：学习并实践自顶向下 (Top-Down) 的分析方法，具体通过**递归下降分析法 (Recursive Descent Parsing)** 来实现。
- 文法转换技巧**：掌握上下文无关文法 (CFG) 的预处理技术，包括消除左递归 (Left Recursion Elimination) 和提取左因子 (Left Factoring)，以满足 LL(1) 文法的要求。
- 程序实现能力**：编写 C 语言程序，能够对包含赋值、循环、条件判断及算术表达式的源代码进行正确的推导。

## b) 内容描述

### 2.1 实验内容概述

本实验旨在实现一个针对类 C 语言子集的语法分析器。该分析器接收词法分析器输出的记号流，依据预定义的 LL(1) 文法规则进行推导，最终输出推导过程以验证语法的正确性。

### 2.2 功能要求

- 输入处理**：读取源代码字符流或 Token 序列。
  - 示例输入：`x = 0.5 * (y + 10 + z) * 3; while(x) { x = x + 1; }`
- 词法与语法协同**：在语法分析过程中动态调用词法分析函数 `getNextToken` 获取当前的 Lookahead 符号。
- 结构分析**：
  - 识别标识符 (ID) 与数字 (Num, 含浮点数)。
  - 解析算术表达式 (优先级与结合性)。
  - 解析控制流语句 (If-Else, While)。
  - 解析复合语句 (代码块 { ... })。
- 输出生成**：打印每一步使用的产生式规则，形成推导序列。

### 2.3 输入输出规范

- 输入**：包含源代码的文本流 (如 ASCII 编码的字符串或文件)。
- 输出**：标准输出显示的推导序列及分析结果。
  - 成功**：按顺序打印产生式，最后提示 "Syntax analysis successful"。
  - 失败**：打印 "Syntax Error" 及具体的错误 Token 信息。

## 2.4 核心文法定义 (CFG)

本实验采用以下经过处理的 LL(1) 文法：

```
1. Program      -> Stmts
2. Stmts        -> Stmt Stmts | epsilon
3. Stmt         -> Block | AssignStmt | IfStmt | WhileStmt
4. Block        -> { Stmts }
5. AssignStmt   -> id = E ;
6. IfStmt       -> if ( E ) Stmt ElsePart
7. ElsePart     -> else Stmt | epsilon
8. WhileStmt    -> while ( E ) Stmt
9. E            -> T E'
10. E'          -> + T E' | epsilon
11. T           -> F T'
12. T'          -> * F T' | epsilon
13. F           -> ( E ) | id | num
```

简要说明：

- **消除左递归**：原数学表达式中的  $E \rightarrow E + T$  被转换为右递归形式  $E \rightarrow T E'$ 。
- **提取左因子**：原 IfStmt 存在  $\text{if}(E) S$  和  $\text{if}(E) S \text{ else } S$  的冲突，通过引入 ElsePart 解决。
- **优先级体现**：通过 E (Expression), T (Term), F (Factor) 的层级自然体现乘法优于加法的优先级。

## c) 设计思路 / 方法

### 3.1 总体设计思路

本实验采用 \*\*递归下降分析法\*\* 实现语法分析器。这是一种自顶向下 (Top-Down) 的分析技术，具有逻辑清晰、易于实现的特点。

1. **非终结符与函数的映射**：文法中的每一个非终结符 (Non-terminal)，如 Program, Stmt, E 等，都对应一个 C 语言函数。
2. **Lookahead 机制**：维护一个全局变量 currentToken 保存当前从词法分析器读取的记号。每个递归函数根据 currentToken 的类型来决定选择哪个产生式。
3. **推导过程**：程序从开始符号 parseProgram() 开始执行，通过函数间的相互递归调用模拟文法的推导过程。
4. **终结符匹配**：定义 match(expectedType) 函数。当产生式右部出现终结符时，调用该函数验证当前 Token 是否匹配，若匹配则读取下一个 Token，否则报错。

### 3.2 文法变换

实验给出的文法已经是经过预处理满足 LL(1) 要求的文法。为了实现递归下降，原始的自然文法 (Natural Grammar) 通常需要进行以下变换：

3.2.1 消除左递归

原算术表达式文法通常包含左递归（如  $E \rightarrow E + T$ ），这会导致递归下降程序陷入循环。

- 原产生式：

```
E -> E + T | T
```

- 变换后：引入尾部非终结符  $E'$ ，改为右递归形式：

```
E -> T E'
E' -> + T E' | epsilon
```

- 同理， $T \rightarrow T * F | F$  也被变换为

```
T -> F T'
T' -> * F T' | epsilon。
```

3.2.2 提取左因子

if 语句通常存在“悬空 else”或前缀冲突问题。

- 原产生式：

```
Stmt -> if (E) Stmt
      Stmt -> if (E) Stmt else Stmt
```

- 冲突点：读取到 if 时，无法决定使用哪条规则。
- 变换后：提取公共前缀 if (E) Stmt，将差异部分放入新非终结符 ElsePart：

```
IfStmt -> if ( E ) Stmt ElsePart
ElsePart -> else Stmt | epsilon
```

3.3 FIRST 集构造

$FIRST(\alpha)$  是指从  $\alpha$  推导出的串中，第一个终结符构成的集合。

根据给出的文法，计算结果如下：

| 非终结符 (x)   | FIRST(x)       | 计算依据                           |
|------------|----------------|--------------------------------|
| F          | { (, id, num } | 直接由产生式 13 得出                   |
| T'         | { *, epsilon } | 由产生式 12 得出                     |
| T          | { (, id, num } | <code>First(T)=First(F)</code> |
| E'         | { +, epsilon } | 由产生式 10 得出                     |
| E          | { (, id, num } | <code>First(E)=First(T)</code> |
| AssignStmt | { id }         | 由产生式 5 得出                      |
| Block      | { { }          | 由产生式 4 得出                      |

| 非终结符 (x)  | FIRST(x)                      | 计算依据  |
|-----------|-------------------------------|---|
| IfStmt    | { if }                        | 由产生式 6 得出                                     |
| WhileStmt | { while }                     | 由产生式 8 得出                                     |
| ElsePart  | { else, epsilon }             | 由产生式 7 得出                                     |
| Stmt      | { id, {, if, while }          | 集合 Assign, Block, If, While 的 First 集         |
| Stmts     | { id, {, if, while, epsilon } | $\text{First}(\text{Stmt}) \cup \{\epsilon\}$ |
| Program   | { id, {, if, while, epsilon } | $\text{First}(\text{Stmts})$                  |

### 3.4 FOLLOW 集构造

FOLLOW(A) 是指在句型中紧跟在非终结符 A 之后的终结符集合。

#### 1. Program:

- 初始加入结束符 \$ (EOF)。
- Follow(Program) = { \$ }**

#### 2. Stmts:

- 作为 Program 的右部：加入 \$。
- 在 Block  $\rightarrow \{ \text{Stmts} \}$  中：加入 }。
- 在 Stmts  $\rightarrow \text{Stmt Stmts}$  中：位于末尾，无新增。
- Follow(Smts) = { }, \$ }**

#### 3. Stmt:

- 在 Stmts  $\rightarrow \text{Stmt Stmts}$  中：加入  $\text{First}(\text{Stmts}) \setminus \{\epsilon\}$ ，即 { id, {, if, while }。
- 因为 Stmts 可推导为空，且 Stmt 是 Stmts 的尾部元素，故加入 Follow(Smts)。
- 在 IfStmt 中位于 else 之前：加入 else。
- Follow(Stmt) = { id, {, if, while, }, \$, else }**

#### 4. ElsePart:

- 位于 IfStmt 末尾，加入 Follow(IfStmt)，即 Follow(Stmt)。
- Follow(ElsePart) = { id, {, if, while, }, \$, else }**

#### 5. E:

- 在 AssignStmt  $\rightarrow \text{id} = \text{E} ;$  中：加入 ;。
- 在  $F \rightarrow ( \text{E} ) / \text{If} / \text{While}$  中：加入 )。
- Follow(E) = { ;, ) }**

#### 6. E':

- 位于 E 末尾，加入 Follow(E)。
- Follow(E') = { ;, ) }**

#### 7. T:

- 在  $E' \rightarrow + T E'$  中: 加入  $\text{First}(E') \setminus \{\epsilon\}$ , 即  $\{+\}$ 。
- 因  $E'$  可为空, 加入  $\text{Follow}(E')$ 。
- $\text{Follow}(T) = \{+, ,, )\}$

8.  $T'$ :

- 位于  $T$  末尾, 加入  $\text{Follow}(T)$ 。
- $\text{Follow}(T') = \{+, ,, )\}$

9.  $F$ :

- 在  $T' \rightarrow * F T'$  中: 加入  $\text{First}(T') \setminus \{\epsilon\}$ , 即  $\{*\}$ 。
- 因  $T'$  可为空, 加入  $\text{Follow}(T')$ 。
- $\text{Follow}(F) = \{*, +, ,, )\}$

| 非终结符     | FOLLOW 集  |
|----------|---|
| Program  | $\{ \$ \}$  |
| Stmts    | $\{ \}, \$ \}$  |
| Stmt     | $\{ \text{id}, \{, \text{if}, \text{while}, \}, \$, \text{else} \}$ |
| ElsePart | $\{ \text{id}, \{, \text{if}, \text{while}, \}, \$, \text{else} \}$ |
| $E$      | $\{ ;, ) \}$  |
| $E'$     | $\{ ;, ) \}$  |
| $T$      | $\{ +, ;, ) \}$   |
| $T'$     | $\{ +, ;, ) \}$   |
| $F$      | $\{ *, +, ;, ) \}$  |

## 3.5 LL(1) 分析表构造

LL(1) 分析表  $M[A,a]$  用于确定当栈顶非终结符为  $A$  且当前输入为  $a$  时, 应选择哪个产生式。构造规则:

- 若  $a \in \text{First}(\alpha)$ , 则  $M[A,a] = A \rightarrow \alpha$ 。
- 若  $\epsilon \in \text{First}(\alpha)$ , 则对所有  $b \in \text{Follow}(A)$ , 令  $M[A,b] = A \rightarrow \alpha$  (即  $A \rightarrow \epsilon$ )。

关键部分的分析表推导:

- 针对 Stmts (产生式 2:  $\text{Stmt Stmts} \mid \epsilon$ ):
  - $\text{First}(\text{Stmt}) = \{ \text{id}, \{, \text{if}, \text{while} \} \rightarrow$  填入  $\text{Stmt Stmts}$ 。
  - $\text{Follow}(\text{Stmts}) = \{ \}, \$ \rightarrow$  填入  $\epsilon$ 。
- 针对  $E'$  (产生式 10:  $+ T E' \mid \epsilon$ ):
  - $\text{First}(+ T E') = \{ + \} \rightarrow$  填入  $+ T E'$ 。
  - $\text{Follow}(E') = \{ ;, ) \} \rightarrow$  填入  $\epsilon$ 。
- 针对 ElsePart (产生式 7:  $\text{else Stmt} \mid \epsilon$ ):

- $\text{First}(\text{else Stmt}) = \{ \text{else} \} \rightarrow$  填入 `else Stmt`。
- $\text{Follow}(\text{ElsePart})$  包含 `{ id, {, if, while, }, $ }`  $\rightarrow$  填入 `epsilon`。
- 注意：这里解决了 `dangling-else` 歧义，遇到 `else` 优先移进，遇到其他符号归约。

## 3.6 递归下降分析法实现

基于上述分析，核心函数的伪代码/C语言实现逻辑如下：

### 3.6.1 语句列表处理

利用 `FIRST` 和 `FOLLOW` 集决定是继续解析下一条语句还是结束当前列表。

```
void parseStmts() {
    // 检查是否属于 First(Stmt)
    if (currentToken.type == TOK_ID || currentToken.type == TOK_LBRACE ||
        currentToken.type == TOK_IF || currentToken.type == TOK_WHILE) {
        // Rule 2: Stmt -> Stmt Stmt
        parseStmt();
        parseStmts();
    }
    // 检查是否属于 Follow(Stmts)
    else if (currentToken.type == TOK_RBRACE || currentToken.type == TOK_EOF) {
        // Rule 2: Stmt -> epsilon
        // Do nothing (consume epsilon)
    }
    else {
        error("Unexpected token in Stmts");
    }
}
```

### 3.6.2 If 语句与 Else 处理

处理 `if` 结构及可选的 `else` 分支。

```
void parseIfStmt() {
    // Rule 6: IfStmt -> if ( E ) Stmt ElsePart
    match(TOK_IF);
    match(TOK_LPAREN);
    parseE();
    match(TOK_RPAREN);
    parseStmt();
    parseElsePart(); // 调用子函数处理 else
}

void parseElsePart() {
    if (currentToken.type == TOK_ELSE) {
        // Rule 7: ElsePart -> else Stmt
        match(TOK_ELSE);
        parseStmt();
    } else {
        // Rule 7: ElsePart -> epsilon
        // 凡是在 Follow(ElsePart) 中的符号，都意味着 else 部分的结束（即为空）
    }
}
```



```

// 这里的 Follow 集包括 }, ;, id, if, while, EOF 等,
// 实际上只要不是 else, 我们都可以认为 ElsePart 为空并返回 (假设没有语法错误)
}
}

```

### 3.6.3 表达式消除左递归处理

处理算术表达式的优先级。

```

void parseE() {
    // Rule 9: E -> T E'
    parseT();
    parseEPrime();
}

void parseEPrime() {
    if (currentToken.type == TOK_PLUS) {
        // Rule 10: E' -> + T E'
        match(TOK_PLUS);
        parseT();
        parseEPrime();
    }
    // Follow(E') = { ), ; }
    else if (currentToken.type == TOK_RPAREN || currentToken.type == TOK_SEMI) {
        // Rule 10: E' -> epsilon
    }
    else {
        // 严格来说, 这里应该报错, 或者让上层调用者处理
    }
}

```

## d) 假设条件

1. **输入格式**: 假设输入是 ASCII 编码的字符流。
2. **词法限制**: 标识符由字母或下划线开头, 数字支持整数和小数, 不支持科学计数法。
3. **语法范围**: 仅支持基本的整型/浮点型算术运算, 不涉及复杂的类型检查或作用域管理。
4. **错误处理**: 采用“恐慌模式”的简化版, 遇到第一个语法错误即报错并终止程序, 不进行复杂的错误恢复。

## e) 相关文法描述

本节详细描述实验所采用的文法定义, 并展示用于构建语法分析器的核心数据——FIRST 集、FOLLOW 集以及最终生成的 LL(1) 预测分析表。

### 1. 完整文法定义

本实验针对的类 C 语言子集文法

$$G = (VN, VT, P, S)$$

定义如下:

- **非终结符 (VN):**

Program, Stmts, Stmt, Block, AssignStmt, IfStmt, ElsePart, WhileStmt, E, E', T, T', F

- **终结符 (VT):**

id, num, =, ;, {, }, if, else, while, (, ), +, \*, \$ (EOF)

- **开始符号 (S):** Program

- **产生式集合 (P):**

1. Program      -> Stmts
2. Stmts       -> Stmt Stmts | epsilon
3. Stmt        -> Block | AssignStmt | IfStmt | WhileStmt
4. Block       -> { Stmts }
5. AssignStmt  -> id = E ;
6. IfStmt       -> if ( E ) Stmt ElsePart
7. ElsePart    -> else Stmt | epsilon
8. WhileStmt   -> while ( E ) Stmt
9. E           -> T E'
10. E'          -> + T E' | epsilon
11. T           -> F T'
12. T'          -> \* F T' | epsilon
13. F           -> ( E ) | id | num

## 2. FIRST 集

FIRST 集决定了当非终结符位于句首时，面对当前的 Lookahead 符号应选择哪条产生式。

| 非终结符       | FIRST 集                       |
|------------|-------------------------------|
| Program    | { id, {, if, while, epsilon } |
| Stmts      | { id, {, if, while, epsilon } |
| Stmt       | { id, {, if, while }          |
| Block      | { { }                         |
| AssignStmt | { id }                        |
| IfStmt     | { if }                        |
| WhileStmt  | { while }                     |
| ElsePart   | { else, epsilon }             |
| E          | { (, id, num }                |
| E'         | { +, epsilon }                |
| T          | { (, id, num }                |
| T'         | { *, epsilon }                |
| F          | { (, id, num }                |

### 3. FOLLOW 集

FOLLOW 集用于处理产生式右部可能推导为空（epsilon）的情况。当非终结符推导为空时，分析器需查看后继符号以决定是否使用空产生式。

| 非终结符       | FOLLOW 集                          |
|------------|-----------------------------------|
| Program    | { \$ }                            |
| Stmts      | { }, \$ }                         |
| Stmt       | { id, {, if, while, }, \$, else } |
| Block      | { id, {, if, while, }, \$, else } |
| AssignStmt | { id, {, if, while, }, \$, else } |
| IfStmt     | { id, {, if, while, }, \$, else } |
| WhileStmt  | { id, {, if, while, }, \$, else } |
| ElsePart   | { id, {, if, while, }, \$, else } |
| E          | { ), ; }                          |
| E'         | { ), ; }                          |
| T          | { +, ), ; }                       |
| T'         | { +, ), ; }                       |
| F          | { *, +, ), ; }                    |

### 4. LL(1) 预测分析表 (Parsing Table)

行表示栈顶的非终结符，列表示当前的 Lookahead 终结符。单元格内容为应当采用的产生式。

（注：空白处表示语法错误；eps 代表 epsilon）

| Non-Term | id            | num | + | * | ( | ) | {             | }   | = | ; | if            | while           | else         | \$    |
|----------|---------------|-----|---|---|---|---|---------------|-----|---|---|---------------|-----------------|--------------|-------|
| Program  | Stmts         |     |   |   |   |   | Stmts         |     |   |   | Stmts         | Stmts           |              | Stmts |
| Stmts    | Stmt<br>Stmts |     |   |   |   |   | Stmt<br>Stmts | eps |   |   | Stmt<br>Stmts | Stmt<br>Stmts   |              | eps   |
| Stmt     | Assign        |     |   |   |   |   | Block         |     |   |   | If            | While           |              |       |
| Block    |               |     |   |   |   |   | { Stmts }     |     |   |   |               |                 |              |       |
| Assign   | id = E ;      |     |   |   |   |   |               |     |   |   |               |                 |              |       |
| IfStmt   |               |     |   |   |   |   |               |     |   |   | if (E)...     |                 |              |       |
| ElsePart | eps           |     |   |   |   |   | eps           | eps |   |   | eps           | eps             | else<br>Stmt | eps   |
| While    |               |     |   |   |   |   |               |     |   |   |               | while<br>(E)... |              |       |

| Non-Term | id   | num  | +      | *      | (     | )   | { | } | = | ;   | if | while | else | \$ |
|----------|------|------|--------|--------|-------|-----|---|---|---|-----|----|-------|------|----|
| E        | T E' | T E' |        |        | T E'  |     |   |   |   |     |    |       |      |    |
| E'       |      |      | + T E' |        |       | eps |   |   |   | eps |    |       |      |    |
| T        | F T' | F T' |        |        | F T'  |     |   |   |   |     |    |       |      |    |
| T'       |      |      | eps    | * F T' |       | eps |   |   |   | eps |    |       |      |    |
| F        | id   | num  |        |        | ( E ) |     |   |   |   |     |    |       |      |    |

关于 ElsePart 和 else 的特别说明：

- 在 M[ElsePart, else] 中，从 First 集看可以使用 ElsePart -> else Stmt，从 Follow 集看可以使用 ElsePart -> epsilon（因为 else 在 Follow 集中）。这是经典的“悬空 else”歧义。
- LL(1) 解决策略：**在此格中优先填入移进动作（else Stmt），即让 else 与最近的 if 匹配。这也是为什么在上表中，else 列对应的是 else Stmt 而不是 eps。

f) 重要数据结构说明

1. 记号类型枚举 (TokenType Enumeration)

用于程序内部唯一标识语法中的每一个终结符。使用枚举代替字符串比较可以显著提高分析器的运行效率，并使 switch-case 逻辑更加清晰。

```
typedef enum {
    /* 关键字 keywords */
    TOK_IF,           // if
    TOK_ELSE,         // else
    TOK_WHILE,        // while

    /* 标识符与字面量 Identifiers & Literals */
    TOK_ID,           // 变量名 (如 x, count)
    TOK_NUM,          // 数值 (如 10, 0.5)

    /* 运算符 Operators */
    TOK_ASSIGN,       // =
    TOK_PLUS,         // +
    TOK_MUL,          // *

    /* 分隔符 Separators */
    TOK_LPAREN,       // (
    TOK_RPAREN,       // )
    TOK_LBRACE,        // {
    TOK_RBRACE,        // }
    TOK_SEMI,         // ;

    /* 控制标记 Control */
}
```

```
TOK_EOF,           // 文件结束 (End of File)
TOK_ERROR          // 词法错误
} TokenType;
```

## 2. 记号结构体 (Token Structure)

这是词法分析器与语法分析器之间的交互单元。它不仅包含记号的类型，还保存了其原始文本值 (lexeme) 和位置信息，以便于后续的语义分析 (如获取数值) 和错误报告。

```
typedef struct {
    TokenType type;           // 记号的类别
    char lexeme[100];        // 记号的实际文本值 (例如 "count", "3.14")
    int lineNo;              // 记号所在的行号 (用于报错定位)
} Token;
```

## 3. 全局分析状态 (Global Parsing State)

为了简化递归函数的参数传递，采用全局变量来维护当前的分析状态。

```
// 当前正在处理的记号 (Lookahead Symbol)
Token currentToken;

// 输入源指针 (可以是文件指针或字符串缓冲区的索引)
FILE *sourceFile;

// 用于词法分析的行号计数器
int currentLine = 1;
```

## g) 核心算法说明

实核心算法基于 **LL(1) 递归下降分析法**。程序的执行流模拟了文法产生式的推导过程，通过“预测”和“匹配”两个动作构建语法结构。

### 1. 终结符匹配算法

这是消耗输入流中 Token 的唯一途径。它验证当前 Token 是否符合预期，如果符合，则命令词法分析器读取下一个 Token；否则抛出语法错误。

```

void match(TokenType expectedType) {
    if (currentToken.type == expectedType) {
        // 1. 打印当前匹配成功的终结符（用于生成推导序列）
        printf("Match: %s\n", currentToken.lexeme);

        // 2. 推进输入流，获取下一个 Token
        getNextToken();
    } else {
        // 3. 错误处理：预期类型与实际类型不符
        syntaxError(expectedType);
    }
}

```

## 2. 语句列表解析算法

该算法展示了如何处理 **递归定义** 和 **空产生式 (Epsilon)**。

- 文法:  $\text{Stmts} \rightarrow \text{Stmt Stmts} \mid \epsilon$
- 逻辑: 查看 `currentToken` 是属于 `FIRST(Stmt)` 还是 `FOLLOW(Stmts)`。

```

void parseStmts() {
    // 检查 currentToken 是否在 FIRST(Stmt) 中
    if (currentToken.type == TOK_ID || currentToken.type == TOK_IF ||
        currentToken.type == TOK_WHILE || currentToken.type == TOK_LBRACE) {

        printf("Stmts -> Stmt Stmts\n"); // 输出推导规则
        parseStmt(); // 处理单条语句
        parseStmts(); // 递归处理剩余语句
    }
    // 检查 currentToken 是否在 FOLLOW(Stmts) 中
    else if (currentToken.type == TOK_RBRACE || currentToken.type == TOK_EOF) {

        printf("Stmts -> epsilon\n"); // 输出推导规则
        // 匹配 epsilon, 不消耗任何 Token, 直接返回
        return;
    }
    else {
        // 既不是语句的开始, 也不是合法的结束符 -> 报错
        error("Unexpected token in Statement List");
    }
}

```

## 3. 表达式解析算法

该算法展示了如何处理 **运算符优先级** 和 **消除左递归后的结构**。

- 文法:  $E \rightarrow T E'$  和  $E' \rightarrow + T E' \mid \epsilon$

```

// 解析表达式 (Expression)
void parseE() {
    printf("E -> T E'\n");
}

```

```

    parseT();          // 先解析优先级较高的项 (Term)
    parseEPrime();     // 再解析低优先级的加法后缀
}

// 解析表达式后缀 (Expression Prime)
void parseEPrime() {
    if (currentToken.type == TOK_PLUS) {
        // 发现加号: E' -> + T E'
        printf("E' -> + T E'\n");
        match(TOK_PLUS); // 消耗 '+'
        parseT();        // 解析右操作数
        parseEPrime();   // 继续寻找后续的加法
    }
    else if (currentToken.type == TOK_RPAREN || currentToken.type == TOK_SEMI) {
        // 发现结束符: E' -> epsilon
        printf("E' -> epsilon\n");
        return;
    }
    else {
        error("Invalid expression syntax");
    }
}
}

```

## 4. 驱动程序

主程序负责初始化并启动分析过程。

1. **初始化**: 打开输入文件，初始化变量。
2. **启动词法**: 调用一次 getNextToken() 获取第一个 Lookahead 符号。
3. **启动语法**: 调用开始符号对应的函数 parseProgram()。
4. **验证结束**: 分析函数返回后，检查 currentToken 是否为 TOK\_EOF，确保整个输入流被正确解析完毕。

## 5. 错误处理

本实验采用简单的 **Panic Mode (恐慌模式)** 的变体：

- 一旦发现语法错误（match 失败或预测表为空），打印包含行号和错误类型的详细信息。
- 停止分析并退出程序（Fail-fast），或者尝试跳过输入直到遇到同步符号（如；或 }）以尝试恢复（可选实现）。

```

void syntaxError(TokenType expected) {
    fprintf(stderr, "Syntax Error at line %d: Expected %s but found '%s'\n",
            currentToken.lineNo,
            tokenTypeToString(expected),
            currentToken.lexeme);
    exit(1); // 终止程序
}

```

## h) 实验核心代码

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <ctype.h>
#include <stdarg.h> // 用于可变参数的日志输出

/* =====
1. 定义与全局变量
===== */

// 定义所有可能的 Token 类型
typedef enum {
    TOK_ID,           // 标识符 (x, y, count)
    TOK_NUM,          // 数字 (10, 0.5)
    TOK_ASSIGN,       // =
    TOK_SEMI,         // ;
    TOK_LPAREN,       // (
    TOK_RPAREN,       // )
    TOK_LBRACE,       // {
    TOK_RBRACE,       // }
    TOK_PLUS,         // +
    TOK_MULT,         // *
    TOK_IF,           // if 关键字
    TOK_ELSE,         // else 关键字
    TOK_WHILE,        // while 关键字
    TOK_EOF,          // 输入结束
    TOK_ERROR         // 词法错误
} TokenType;

// Token 结构体
typedef struct {
    TokenType type;
    char str[100]; // 存储 Token 的原始字符串值
} Token;

// 全局变量
char inputBuffer[2048]; // 输入字符流缓冲区
int pos = 0;           // 当前字符流读取位置
Token currentToken;     // 当前正在分析的 Token
FILE *fileOut = NULL;   // 输出文件指针

/* =====
2. 函数声明
===== */

// 工具函数
void logOutput(const char *format, ...); // 同时输出到屏幕和文件
void getNextToken();                    // 词法分析器 (Scanner)
void match(TokenType expected);         // 匹配终结符
void error(const char *msg);           // 报错处理
```



```

// 语法分析函数（对应每个非终结符）
void parseProgram();    // 起始符号
void parseStmt();       // 语句（Statement）
void parseBlock();      // 代码块（Block）
void parseStmts();      // 语句列表
void parseAssignStmt(); // 赋值语句
void parseIfStmt();     // 条件语句
void parseElsePart();   // else 部分（提取左因子后）
void parsewhileStmt();  // 循环语句

// 表达式分析（消除左递归后的文法）
void parseE(); // E -> T E'
void parseE_(); // E' -> + T E' | epsilon
void parseT(); // T -> F T'
void parseT_(); // T' -> * F T' | epsilon
void parseF(); // F -> ( E ) | id | num

/* =====
3. 主函数（入口）
===== */
int main() {
    // 1. 打开输出文件
    fileOut = fopen("output.txt", "w");
    if (fileOut == NULL) {
        printf("Error: Could not create output.txt\n");
        return 1;
    }

    // 2. 准备输入数据（模拟字符流）
    // 你可以在这里修改测试用例
    // 测试用例涵盖：赋值，算术运算，括号，循环结构
    strcpy(inputBuffer, "x = 0.5 * (y + 10 + z) * 3; while(x) { x = x + 1; }");

    logOutput("Lab2 LL(1) Syntax Parser\n");
    logOutput("-----\n");
    logOutput("Input Stream: %s\n", inputBuffer);
    logOutput("-----\n");
    logOutput("Sequence of Derivations (Syntax Tree Construction):\n\n");

    // 3. 初始化并启动分析
    pos = 0;
    getNextToken(); // 预读第一个 Token
    parseProgram(); // 进入递归下降分析

    // 4. 结束检查
    if (currentToken.type == TOK_EOF) {
        logOutput("\n-----\n");
        logOutput("Result: SUCCESS. Syntax tree built successfully.\n");
        logOutput("Output saved to 'output.txt'.\n");
    } else {
        error("Unexpected characters at the end of input.");
    }
}

```

```

    fclose(fileOut);
    return 0;
}

/* =====
4. 工具函数实现
===== */

// 双路输出日志函数（屏幕 + 文件）
void logOutput(const char *format, ...) {
    va_list args;

    // 输出到控制台
    va_start(args, format);
    vprintf(format, args);
    va_end(args);

    // 输出到文件
    if (fileOut != NULL) {
        va_start(args, format);
        vfprintf(fileOut, format, args);
        va_end(args);
    }
}

// 报错函数
void error(const char *msg) {
    logOutput("\n[Syntax Error] %s (Current Token: '%s')\n", msg, currentToken.str);
    if (fileOut) fclose(fileOut);
    exit(1);
}

// 终结符匹配函数
void match(TokenType expected) {
    if (currentToken.type == expected) {
        getNextToken(); // 匹配成功，读取下一个 Token
    } else {
        char errBuf[100];
        sprintf(errBuf, "Expected token type %d but found '%s'", expected,
currentToken.str);
        error(errBuf);
    }
}

/* =====
5. 词法分析器 (Lexer)
===== */
void getNextToken() {
    // 1. 跳过空白字符（空格，Tab，换行）
    while (inputBuffer[pos] == ' ' || inputBuffer[pos] == '\t' ||
        inputBuffer[pos] == '\n' || inputBuffer[pos] == '\r') {
        pos++;
    }
}

```

```

}

// 2. 判断字符串结束
if (inputBuffer[pos] == '\0') {
    currentToken.type = TOK_EOF;
    strcpy(currentToken.str, "EOF");
    return;
}

char c = inputBuffer[pos];

// 3. 识别标识符 (Identifier) 或 关键字 (Keywords)
if (isalpha(c) || c == '_') {
    int i = 0;
    while (isalnum(inputBuffer[pos]) || inputBuffer[pos] == '_') {
        currentToken.str[i++] = inputBuffer[pos++];
    }
    currentToken.str[i] = '\0';

    // 检查是否是保留字
    if (strcmp(currentToken.str, "if") == 0) currentToken.type = TOK_IF;
    else if (strcmp(currentToken.str, "else") == 0) currentToken.type = TOK_ELSE;
    else if (strcmp(currentToken.str, "while") == 0) currentToken.type = TOK_WHILE;
    else currentToken.type = TOK_ID;
    return;
}

// 4. 识别数字 (Numbers, 支持小数)
if (isdigit(c)) {
    int i = 0;
    while (isdigit(inputBuffer[pos]) || inputBuffer[pos] == '.') {
        currentToken.str[i++] = inputBuffer[pos++];
    }
    currentToken.str[i] = '\0';
    currentToken.type = TOK_NUM;
    return;
}

// 5. 识别运算符和界符
// 每次处理完后 pos++ 移动指针
currentToken.str[0] = c;
currentToken.str[1] = '\0';
pos++;

switch(c) {
    case '=': currentToken.type = TOK_ASSIGN; break;
    case ';': currentToken.type = TOK_SEMI; break;
    case '+': currentToken.type = TOK_PLUS; break;
    case '*': currentToken.type = TOK_MULT; break;
    case '(': currentToken.type = TOK_LPAREN; break;
    case ')': currentToken.type = TOK_RPAREN; break;
    case '{': currentToken.type = TOK_LBRACE; break;
    case '}': currentToken.type = TOK_RBRACE; break;

```

```

        default: currentToken.type = TOK_ERROR; break;
    }
}

/* =====
6. 语法分析器 (LL(1) Logic)
===== */

// Grammar: Program -> Stmts
void parseProgram() {
    logOutput("Program -> Stmts\n");
    parseStmts(); // 解析语句列表
}

// Grammar: Stmts -> Stmt Stmts | epsilon
// 解释: 如果当前 token 属于 First(Stmt), 则解析 Stmt; 否则如果是 '}' 或 EOF, 则推导为空。
void parseStmts() {
    // First(Stmt) = { id, if, while, { }
    if (currentToken.type == TOK_ID || currentToken.type == TOK_IF ||
        currentToken.type == TOK_WHILE || currentToken.type == TOK_LBRACE) {
        logOutput("Stmts -> Stmt Stmts\n");
        parseStmt();
        parseStmts();
    } else {
        // Follow(Stmts) = { '}', EOF }
        logOutput("Stmts -> epsilon\n");
    }
}

// Grammar: Stmt -> Block | AssignStmt | IfStmt | WhileStmt
void parseStmt() {
    if (currentToken.type == TOK_LBRACE) {
        logOutput("Stmt -> Block\n");
        parseBlock();
    }
    else if (currentToken.type == TOK_ID) {
        logOutput("Stmt -> AssignStmt\n");
        parseAssignStmt();
    }
    else if (currentToken.type == TOK_IF) {
        logOutput("Stmt -> IfStmt\n");
        parseIfStmt();
    }
    else if (currentToken.type == TOK_WHILE) {
        logOutput("Stmt -> WhileStmt\n");
        parseWhileStmt();
    }
    else {
        error("Expected start of a statement (id, if, while, {)");
    }
}

// Grammar: Block -> { Stmts }

```

```

void parseBlock() {
    logOutput("Block -> { Stmts }\n");
    match(TOK_LBRACE);
    parseStmts();
    match(TOK_RBRACE);
}

// Grammar: AssignStmt -> id = E ;
void parseAssignStmt() {
    logOutput("AssignStmt -> id = E ;\n");
    match(TOK_ID);      // 匹配 id
    match(TOK_ASSIGN);  // 匹配 =
    parseE();           // 解析表达式
    match(TOK_SEMI);    // 匹配 ;
}

// Grammar: WhileStmt -> while ( E ) Stmt
void parsewhileStmt() {
    logOutput("whilestmt -> while ( E ) Stmt\n");
    match(TOK_WHILE);
    match(TOK_LPAREN);
    parseE();
    match(TOK_RPAREN);
    parseStmt(); // 循环体
}

// Grammar: IfStmt -> if ( E ) Stmt ElsePart
// 注意：提取左因子，处理 else
void parseIfStmt() {
    logOutput("IfStmt -> if ( E ) Stmt ElsePart\n");
    match(TOK_IF);
    match(TOK_LPAREN);
    parseE();
    match(TOK_RPAREN);
    parseStmt(); // if 的主体
    parseElsePart(); // 处理可能的 else
}

// Grammar: ElsePart -> else Stmt | epsilon
void parseElsePart() {
    if (currentToken.type == TOK_ELSE) {
        logOutput("ElsePart -> else Stmt\n");
        match(TOK_ELSE);
        parseStmt();
    } else {
        // 如果不是 else, 推导为空 (epsilon)
        logOutput("ElsePart -> epsilon\n");
    }
}

/* --- 算术表达式处理 (消除左递归) --- */

// Grammar: E -> T E'

```

```

void parseE() {
    logOutput("E -> T E'\n");
    parseT();
    parseE_();
}

// Grammar: E' -> + T E' | epsilon
void parseE_() {
    if (currentToken.type == TOK_PLUS) {
        logOutput("E' -> + T E'\n");
        match(TOK_PLUS);
        parseT();
        parseE_();
    } else {
        // Follow(E') 包括 ), ;
        logOutput("E' -> epsilon\n");
    }
}

// Grammar: T -> F T'
void parseT() {
    logOutput("T -> F T'\n");
    parseF();
    parseT_();
}

// Grammar: T' -> * F T' | epsilon
void parseT_() {
    if (currentToken.type == TOK_MULT) {
        logOutput("T' -> * F T'\n");
        match(TOK_MULT);
        parseF();
        parseT_();
    } else {
        // Follow(T') 包括 +, ), ;
        logOutput("T' -> epsilon\n");
    }
}

// Grammar: F -> ( E ) | id | num
void parseF() {
    if (currentToken.type == TOK_LPAREN) {
        logOutput("F -> ( E )\n");
        match(TOK_LPAREN);
        parseE();
        match(TOK_RPAREN);
    } else if (currentToken.type == TOK_ID) {
        logOutput("F -> id (%s)\n", currentToken.str);
        match(TOK_ID);
    } else if (currentToken.type == TOK_NUM) {
        logOutput("F -> num (%s)\n", currentToken.str);
        match(TOK_NUM);
    } else {

```

```
        error("Invalid Factor. Expected '(', identifier, or number.");
    }
}
```


## i) 运行测试用例

输入:

```
x = 0.5 * (y + 10 + z) * 3; while(x) { x = x + 1; }
```

程序运行输出 (output.txt):

```
Syntax_Parser > ≡ output.txt
1  Lab2 LL(1) Syntax Parser
2  -----
3  Input Stream: x = 0.5 * (y + 10 + z) * 3; while(x) { x = x + 1; }
4  -----
5  Sequence of Derivations (Syntax Tree Construction):
6
7  Program -> Stmts
8  Stmts -> Stmt Stmts
9  Stmt -> AssignStmt
10 AssignStmt -> id = E ;
11 E -> T E'
12 T -> F T'
13 F -> num (0.5)
14 T' -> * F T'
15 F -> ( E )
16 E -> T E'
17 T -> F T'
18 F -> id (y)
19 T' -> epsilon
20 E' -> + T E'
21 T -> F T'
22 F -> num (10)
23 T' -> epsilon
24 E' -> + T E'
25 T -> F T'
26 F -> id (z)
27 T' -> epsilon
28 E' -> epsilon
```

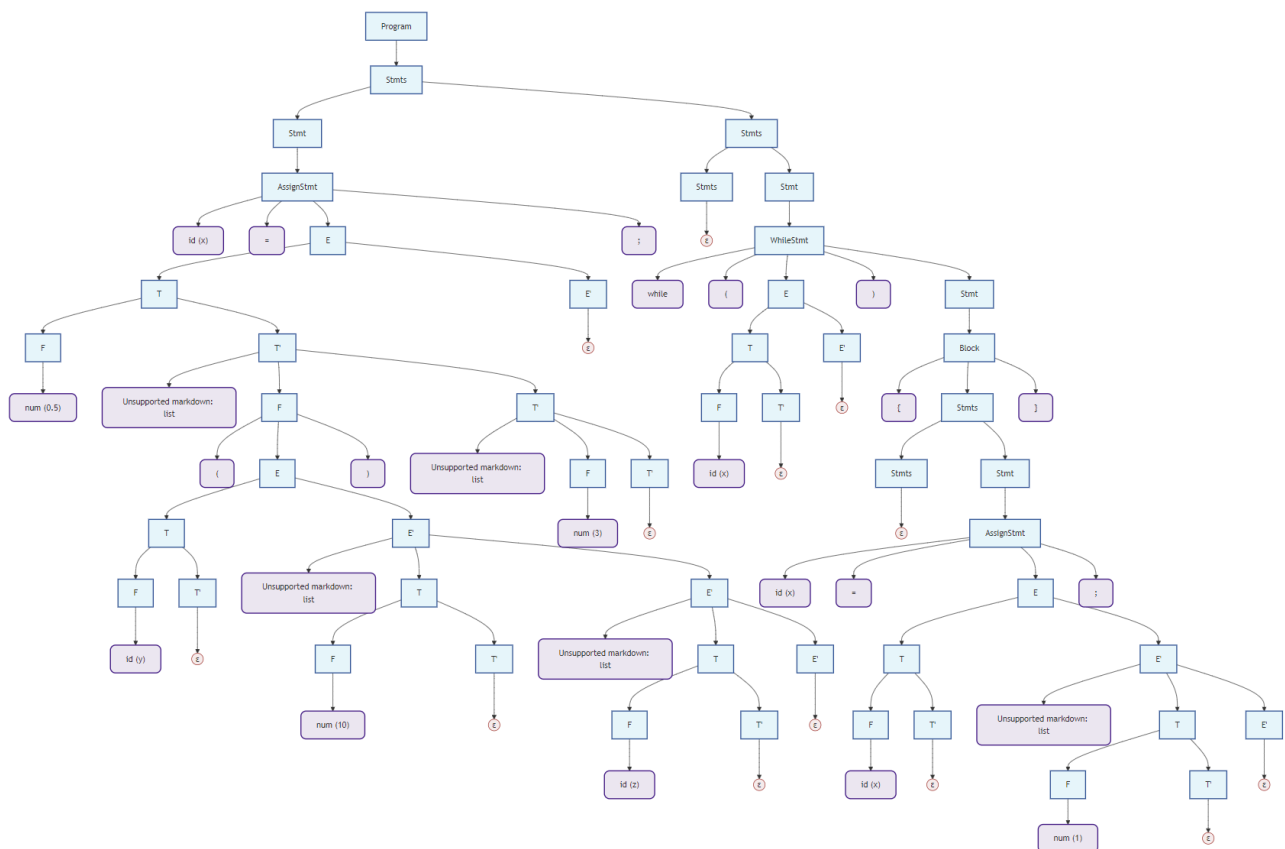
Syntax\_Parser >  output.txt

```

29  T' -> * F T'
30  F -> num (3)
31  T' -> epsilon
32  E' -> epsilon
33  Stmts -> Stmt Stmts
34  Stmt -> WhileStmt
35  WhileStmt -> while ( E ) Stmt
36  E -> T E'
37  T -> F T'
38  F -> id (x)
39  T' -> epsilon
40  E' -> epsilon
41  Stmt -> Block
42  Block -> { Stmts }
43  Stmts -> Stmt Stmts
44  Stmt -> AssignStmt
45  AssignStmt -> id = E ;
46  E -> T E'
47  T -> F T'
48  F -> id (x)
49  T' -> epsilon
50  E' -> + T E'
51  T -> F T'
52  F -> num (1)
53  T' -> epsilon
54  E' -> epsilon
55  Stmts -> epsilon
56  Stmts -> epsilon

```

### 语法分析树 (parsing tree) :



### 结果分析:



一. 分析器成功识别了 Program 由两条主要语句组成：

1. **赋值语句 (AssignStmt)**:  $x = \dots$  ;
2. **循环语句 (WhileStmt)**:  $\text{while}(x) \{ \dots \}$

这体现在推导序列第 8 行  $\text{Stmts} \rightarrow \text{Stmt Stmts}$  (解析第一条语句) 和第 33 行  $\text{Stmts} \rightarrow \text{Stmt Stmts}$  (解析第二条语句) 的递归调用上, 最终在第 56 行以  $\text{Stmts} \rightarrow \epsilon$  结束, 表明整个输入流符合文法结构。

二. 在处理表达式  $0.5 * (y + 10 + z) * 3$  时, 分析器正确体现了 LL(1) 文法的层级设计:

- **括号优先级**: 在第 15 行  $F \rightarrow (E)$ , 分析器在处理乘法  $*$  的过程中遇到括号, 强制优先进入  $E$  (Expression) 层级进行加法运算, 正确实现了括号改变优先级的逻辑。
- **连续运算**: 对于  $y + 10 + z$ , 通过  $E' \rightarrow + T E'$  的右递归结构 (第 20 行和第 24 行) 成功解析了连续相加, 验证了消除左递归后的文法能正确处理由多个同级运算符组成的链式表达式。

三. 复合语句与嵌套

- **While 循环解析**: 程序正确识别了 while 关键字、括号内的条件表达式  $E$  以及循环体  $\text{Stmt}$ 。
- **代码块 (Block) 处理**: 在第 41-42 行,  $\text{Stmt} \rightarrow \text{Block}$  和  $\text{Block} \rightarrow \{ \text{Stmts} \}$  的推导表明分析器能够处理嵌套的作用域。循环体内的  $x = x + 1$  被作为一个新的  $\text{Stmts}$  序列解析, 证明了文法支持语句的无限嵌套 (例如 while 中还可以包含 if 或另一个 while)。

## j) 遇到的问题及解决方案

### 1. 问题: 左递归导致的无限循环

- 描述: 最初按照数学习惯定义文法  $E \rightarrow E + T$ , 导致  $\text{parseE}()$  刚开始就无限调用  $\text{parseE}()$ , 造成栈溢出。
- 解决: 将文法改写为右递归形式  $E \rightarrow T E'$ , 并引入空产生式  $\epsilon$ 。

### 2. 问题: 悬空 Else (Dangling Else)

- 描述: 解析  $\text{if}(E) \text{ if}(E) S \text{ else } S$  时, 无法确定  $\text{else}$  属于哪个  $\text{if}$ 。
- 解决: 通过提取左因子, 将  $\text{if}$  语句定义为必须包含  $\text{ElsePart}$ 。在代码中,  $\text{parseElsePart}$  优先匹配  $\text{else}$  关键字, 这隐式地实现了“ $\text{else}$  与最近的未匹配  $\text{if}$  匹配”的规则 (贪心匹配)。

### 3. 问题: 输出格式混乱

- 描述: 需要同时在控制台查看进度, 并保存结果到文件, 单纯用  $\text{printf}$  不够。
- 解决: 封装了  $\text{logOutput}$  函数, 使用  $\text{stdarg.h}$  库处理可变参数, 实现了一次调用双重输出。

## k) 心得体会

通过本次实验, 我从理论到实践完整地走通了语法分析的流程。

1. **文法即代码**: 我深刻体会到了递归下降分析法的优美之处——文法的结构直接映射为代码的函数调用结构。只要文法设计得当 (满足 LL(1)), 代码编写几乎是机械式的翻译过程。
2. **预处理的重要性**: 实验中最困难的部分不是写代码, 而是设计文法。消除左递归和二义性是保证分析器正常工作的前提。
3. **局限性**: LL(1) 文法对写法的限制较大 (不能左递归), 这使得表达式的文法变得不如原始文法直观 (如  $E'$  的引入)。相比之下, LR 分析法虽然手工构造复杂, 但能处理更广泛的文法, 这是后续值得深入学习的方向。