

哈尔滨工业大学（深圳）

**《编译原理》  
实验指导书**

2024 秋

## 目录

前言 .....	2
第一部分 .....	3
一 实验整体概述 .....	3
1.1 目录说明 .....	3
1.2 代码架构示意图 .....	5
1.3 源语言 .....	5
1.4 目标平台 && 目标语言 .....	6
1.5 关于 NotImplementedException .....	6
二 实验环境及代码规范 .....	7
2.1 实验环境 .....	7
2.2 Java 语言标准 .....	8
2.3 Java 代码注意事项 .....	8
2.3 Java 程序设计风格 .....	9
三 脚本工具 .....	10
3.1 快速比对脚本 .....	10
四 问题反馈 .....	11
第二部分 .....	11
实验一 词法分析器的实现 .....	11
1.1 实验题目 .....	11
1.2 实验目的 .....	11
1.3 实验内容 .....	12
1.4 分析与设计 .....	12
1. 什么是关键字（也是保留字）？ .....	12
2. 什么是标识符？ .....	12
3. 什么是常数？ .....	12
4. 什么是运算符？ .....	13
5. 什么是界符？ .....	13
1.5 实验总体步骤 .....	13
1.5.1 创建编码表 .....	13
1.5.2 创建类 C 语言文法 .....	13
1.5.3 有限自动机 .....	14
1.6 实验框架概述 .....	16
1.6.1 编译程序整体框架图 .....	16
1.6.2 词法分析程序的输入 .....	17
1.6.3 词法分析程序的输出 .....	18
1.7 实验框架设计简介 .....	20
1.7.1 数据结构说明 .....	20
1.7.2 主程序流程说明 .....	24
1.8 实验框架代码说明 .....	26
1.8.1 词法分析程序框架代码 UML 图 .....	26
1.8.2 主要类说明 .....	27
1.8.3 补充说明 .....	29

## 前言

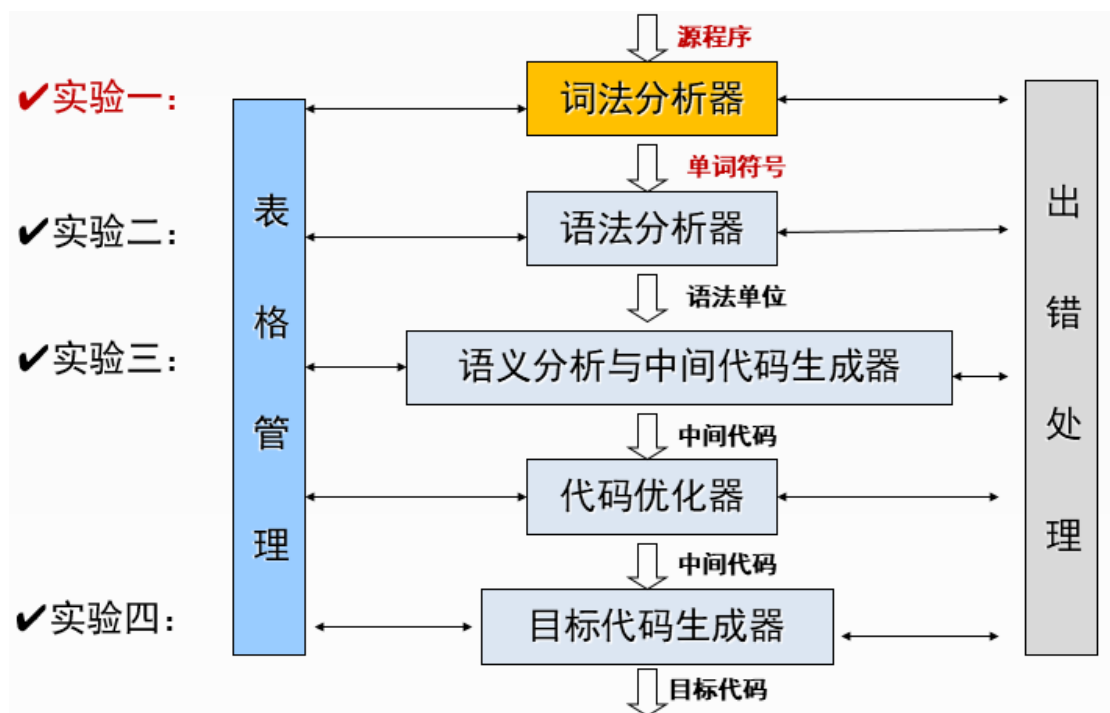
本学期的编译原理实验总共有 4 次实验，全部为设计型实验。

实验一：词法分析器的实现；

实验二：自底向上的语法分析—LR（1）；

实验三：典型语句的语义分析及中间代码生成；

实验四：目标代码生成。



## 第一部分

### 一 实验整体概述

一个 Java 实现的 TXTv2 语言编译器，目标平台为 RISC-V 32（指令集 RV32M）。

#### 1.1 目录说明

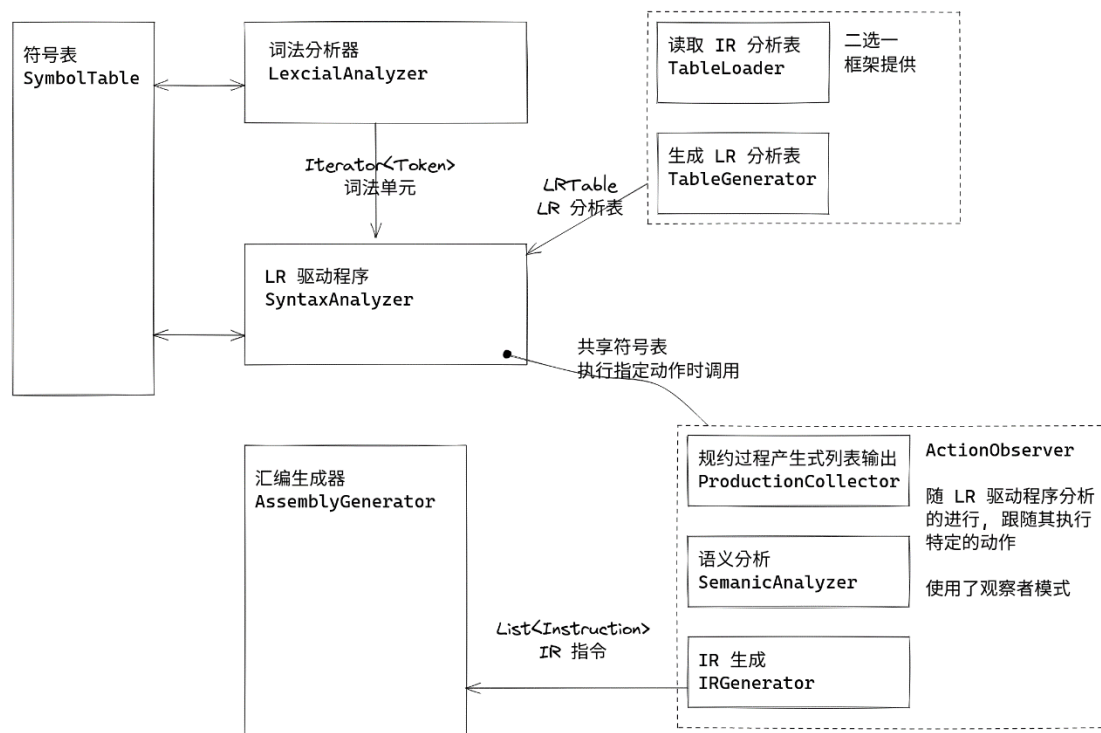
```
tree . -L 2 --dirsfirst --sort=name
```

```
.
├── data
│   ├── in          # 提供给程序的输入数据
│   ├── out         # 程序的输出数据
│   └── std          # 用作参考的标准输出数据
├── scripts
│   ├── check-result.py # 对输出进行 check 的脚本
│   ├── diff.py        # 忽略文件尾空行和行首位空白符的 diff 工具
│   └── make-template.py # 从代码出框架的脚本, 无需关注
└── src               # 源码目录
```

```
tree src/cn/edu/hitsz/compiler -L 2 --dirsfirst --sort=name
```

src/cn/edu/hitsz/compiler	
├── asm	# 实验四: 汇编生成
│   └── AssemblyGenerator.java	# 汇编生成器
├── ir	# 基础架构: IR
│   ├── BasicBlock.java	# IR 中的基本块
│   ├── Instruction.java	# IR 中的指令
│   ├── InstructionKind.java	# IR 中的指令类型
│   ├── IRImmediate.java	# IR 中的立即数
│   ├── IRValue.java	# IR 中的值
│   ├── IRVariable.java	# IR 中的变量
│   └── Program.java	# IR 中代表完整的一个程序的类
├── lexer	# 实验一: 词法分析
│   ├── LexicalAnalyzer.java	# 词法分析器
│   ├── Token.java	# 词法单元
│   └── TokenKind.java	# 词法单元类别
├── parser	# 实验二/三: 语法分析, 语义分析, IR 生成
│   ├── table	# 读取/生成 LR 表的工具
│   ├── ActionObserver.java	# 观察者接口
│   ├── IRGenerator.java	# IR 生成
│   ├── ProductionCollector.java	# 规约产生式记录
│   ├── SemanticAnalyzer.java	# 语义分析
│   └── SyntaxAnalyzer.java	# 语法分析
├── pass	# 可选部分: 优化 pass
│   └── AdjustIRForRISCV.java	# 实验四 (可选部分): 体系结构相关优化与
IR 调整	
│   ├── IRPass.java	# 所有 IR 优化 pass 的基类
│   └── IRPassManager.java	# IR 优化 pass 管理器
├── symtab	# 符号表
│   ├── SourceCodeType.java	# 源语言变量类型
│   ├── SymbolTableEntry.java	# 符号表条目
│   └── SymbolTable.java	# 符号表
├── utils	# 杂项/工具
│   └── CasePaths.java	# 记录程序中用到的相对于测试用例的各
类路径	
│   ├── FileUtils.java	# 文件读写工具
│   ├── IREmulator.java	# 评测用 IR 解释器
│   └── ListForwarder.java	# 辅助类, 用于便捷地为某个类实现 List
接口	
├── Main.java	# 主函数
├── MainForAutoJudge.java	# 用于自动化评测的主函数
└── NotImplementedException.java	# 用于填充待实现部分的异常

## 1.2 代码架构示意图



## 1.3 源语言

由于实验种种限制，本次实验实现的简单编译器的源语言支持的语法包含变量的申明、赋值、简单算术表达式这些简单语法。

```

int result;
int a;
int b;
int c;
a = 8;
b = 5;
c = 3 - a;
result = a * b - (3 + b) * (c - a);
return result;
  
```

## 1.4 目标平台 && 目标语言

我们的目标是生成在支持 RV32M 指令集的 RISC-V 32 机器上可以成功运行的汇编代码。然而这份代码要成功运行还得通过汇编器和链接器生成可执行文件并链接到相应的库上。

然而在实际操作中我们并没有一台 RISC-V 32 的机器，就算通过 qemu 模拟，我们仍不可避免地需要与 GNU 汇编伪指令操作打交道。为了避免这些麻烦，我们使用 RARS(RISC-V Assembler and Runtime Simulator) 进行模拟。可以认为我们的目标机器为 RARS。

根据 RISC-V 的寄存器使用约定，我们使用 Caller 保存的 t0 到 t6，亦即 x5-7, x28-31 作为汇编代码中任意使用的寄存器，使用 a0，亦即 x10 作为程序的返回值。并且我们使用 CompactDataAtZero 选项规定 data 的开始地址为 0x0。通过 RARS 我们可以这样进行汇编程序的执行：

```
$ java -jar rars.jar mc CompactDataAtZero a0 nc dec ae255 riscv1.asm
Program terminated by dropping off the bottom.
a0 10945
```

## 1.5 关于 NotImplementedException

这个 NotImplementedException 的异常在期待你实现但是你还没有实现的函数中被抛出。

由于我们一次性下发了整个项目的代码，所以在你完成整个实验之前，程序都会以抛出 NotImplementedException 异常然后结束。

注意：在完成实验一代码后，因为后续实验二、三、四在主函数中已调用但具体实现代码并未完成，所以依然会在后续代码中抛出 NotImplementedException 的异常，请你不必介意。如若不想让未完成的实验代码抛出异常，那你可以用以下方式修改主函数代码：

```
public class Main {  
    public static void main(String[] args) {  
        // 构建符号表以供各部分使用  
        TokenKind.loadTokenKinds();  
        final var symbolTable = new SymbolTable();  
  
        // 词法分析  
        final var lexer = new LexicalAnalyzer(symbolTable);  
        lexer.loadFile(FilePathConfig.SRC_CODE_PATH);  
        lexer.run();  
        lexer.dumpTokens(FilePathConfig.TOKEN_PATH);  
        final var tokens = lexer.getTokens();  
        symbolTable.dumpTable(FilePathConfig.SYMBOL_TABLE_PATH);  
  
        /**因为后续实验代码未完成，在主函数中调用会抛 NonImplementedException 异常，你可以选择先注释后续调用，解决抛异常的问题*/  
  
        /**  
        //实验二  
        //实验三  
        //实验四  
        */  
    }  
}
```

如果你完成了每次实验的要求，那么这种结束方式是无关大雅的。倘若你担心遗漏了某次实验中待填充的函数，那么请你先清空 `data/out` 文件夹，随后运行程序，再检查一下程序运行之后在 `data/out` 中的输出情况。只要输出的文件数量符合要求且数量/内容是正确的，那么该实验就算是完成了。

## 二 实验环境及代码规范

### 2.1 实验环境

1. Java 17 及以上版本.
2. RARS.
3. 编译工作台 (可选)



## 2.2 Java 语言标准

使用 Java 最新的 LTS 版本，即 Java 17。Java 17 以上以及 Java 17 处于 preview 的特性都不建议使用，方便统一批改。

鼓励使用 Java 的新特性以及新的标准 API，如：

- ✧ 模式匹配 instanceof (Java 16)
- ✧ 增强的 switch 表达式 (Java 14)
- ✧ record 类(Java 16)
- ✧ Stream API 对元素流进行函数式操作 (Java 8)

## 2.3 Java 代码注意事项

Java 是一个很 OOP 的语言，在使用 Java 完成本实验时请注意之前所学的面向对象六大原则。

### 1. 依赖倒转原则(DIP)

依赖倒转原则(DIP)是指一种特定的解耦(传统的依赖关系建立在高层次上,而具体的策略设置则应用在低层次的模块上)形式，使得高层次的模块不依赖于低层次的模块的实现细节，依赖关系被颠倒(反转)，从而使得低层次模块依赖于高层次模块的需求抽象。

高层次的模块不应该依赖于低层次的模块，两者都应该依赖于抽象接口。抽象接口不应该依赖于具体实现，而具体实现则应该依赖于抽象接口。

#### 好的实现

```
List<Instruction> porcessInstruction(List<Instruction> originInstructions){  
    // some code  
    return res;  
}
```

#### 坏的实现

```
ArrayList<Instruction> porcessInstruction(ArrayList<Instruction> originInstructions){  
    // some code  
    return res;  
}
```

## 2. 底米特原则

- ✧ 每个对象应该对其他对象尽可能最少的知道
- ✧ 每个对象应该仅和其朋友通信；不和陌生人通信
- ✧ 仅仅和直接朋友通信

好的实现

将各个模块统一交由顶层模块进行调用，各个模块通过顶层模块使用约定好的 API 接口进行传参。

坏的实现

各个模块相互调用，如链式调用：当 A、B、C 为同一层级的模块时，Main call A、A call B、B call C ...

## 2.3 Java 程序设计风格

Java 是 Java。换句话讲，Java 不是 C。

没必要什么数据结构都用数组(整数索引的表)来实现，特别是当这个"表" 实际上具有更深层的结构的时候。

深入地，没必要什么东西都用一个整数来代表。如果一个对象具有任何与整数不同的性质，那它就不是整数。尽量使用对象的引用来引用对象，而不是拿着一个 int 再去什么地方查询。

字符串是方便使用且默认不可变的，使用字符串与使用任何类型一样方便，字符串不再是危险且需要托管的。字符(Character)跟字节(Byte)是不同的，前者是任何在特定编码下代表某种自然语言中的最小单元的数据，后者是一个没有附加结构的 8bit 数据。不要假设所有字符都是 ASCII 编码下的合法字符。

通用数据结构(容器/集合类)是方便且随时可用的，应该尽可能根据职责拆分需要通用数据结构维护的数据，而不是集中写一个数据结构来维护所有信息。

对外隐藏是方便必要的，每个标识符的访问范围都应该尽可能地小。

特殊值是危险的，在任何需要暴露到包外的接口中，请使用 Optional 代替返回 null/-1，或者使用异常代替错误值。

对于本实验:

- ✧ 性能大部分时候是无关紧要的, 请尽可能优雅地实现;
- ✧ 具体到本实验而言, 输入保证所有字符都是合法的 ASCII 编码。你应尽可能使用 `Character` 中的方法, 而非直接判断 ASCII 编码大小。

## 三 脚本工具

### 3.1 快速比对脚本

在完成实验时, 我们可能会反复执行程序, 输出结果, 然后比对标准输出。然而人眼比对文本是一个较为繁琐的过程, 我们在 `scripts` 文件夹下提供了两个快速比对的工具。

#### **diff.py**

Usage:

```
python -u diff.py /path/for/stdfile /path/for/srcfile
```

比对 `/path/for/stdfile` 与 `/path/for/srcfile` 的差异, 忽略文件尾空行和行首尾空白符, 若有差异则输出首行不匹配。

#### **check-result.py**

Usage:

```
python -u check-result.py n /path/for/stdfiles /path/for/outputfiles
```

其中 `n` 为实验 `n` 的序号, 如 4 代表实验四: `/path/for/stdfiles` 与 `/path/for/outputfiles` 是标准文件的文件夹和输出文件的文件夹。

此脚本会比对本次实验与之前所有实验的输出文件和标准文件的差异。特别地, 对于实验四, 我们直接输出了 `rars` 的输出, 需要自行与程序的预期返回值比对。

脚本固定了 `rars.jar` 的路径, 需要修改为你电脑上 `rars.jar` 的路径:

`rars_path = "/home/test/rars.jar"`

## 四 问题反馈

实验框架从一份已经过正确性检验的代码通过删改生成，框架的正确性已经过初步检验，注释部分可能存在部分错别字或措辞不当。

实验指导书可能存在部分错别字或措辞不准确，可能存在未说明清楚的部分，或某些部分如给出的示例实现部分未经过正确性验证，可能存在错误。若发现实验框架或实验指导书中存在任何错误，或你对其有任何疑问，欢迎提交 issue 或直接向老师或助教反馈，我们将尽快解决。

## 第二部分

### 实验一 词法分析器的实现

#### 1.1 实验题目

手工设计类 C 语言的词法分析器（可以是 c 语言的子集）。

#### 1.2 实验目的

1. 加深对词法分析程序的功能及实现方法的理解；
2. 对类 C 语言的文法描述有更深入的认识，理解有穷自动机、编码表和符号表在编译的整个过程中的应用；
3. 设计并编程实现一个词法分析程序，对类 C 语言源程序段进行词法分析，加深对高级语言的认识。
4. 实验学时数：2 学时。

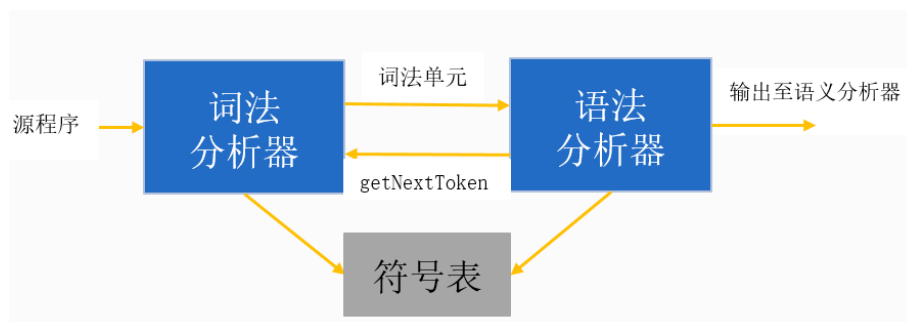
### 1.3 实验内容

编写一个词法分析程序，读取代码文件，对文件内自定义的类 C 语言程序段进行词法分析。处理 C 语言源程序，过滤掉无用符号，分解出正确的单词，以二元组形式存输出放在文件中。

1. 词法分析程序输入：以文件形式存放自定义的类 C 语言程序段；
2. 词法分析程序输出：以文件形式存放的 Token 串和简单符号表；
3. 词法分析程序输入单词类型要求：输入的 C 语言程序段包含常见的关键字，标识符，常数，运算符和分界符等。

### 1.4 分析与设计

要手工设计词法分析器，实现 C 语言子集的认识，就要明白什么是词法分析器，它的功能是什么。



词法分析是编译程序进行编译时第一个要进行的任务，主要是对源程序进行编译预处理（去除注释、无用的回车换行找到包含的文件等）之后，对整个源程序进行分解，分解成一个个单词。这些单词只有五类，分别是标识符、保留字、常数、运算符、界符。可以说词法分析面向的对象是单个字符，目的是把它们组成有效的单词（字符串），从而作为语法分析的输入来分析是否符合语法规则，并且进行语法制导的语义分析产生中间代码，进而优化并生成目标代码。

综上：词法分析器的输入是源代码字符流，输出是单词序列。

1. 什么是关键字（也是保留字）？

关键字是又程序语言定义的具有固定意义的标识符。例如 C 语言中的 if, else, for, while 都是保留字，这些字通常不用作一般标识符。

2. 什么是标识符？

标识符用来表示各种名字，如变量名，数组名，过程名等。

3. 什么是常数？

常数的类型一般是有整型，实型、布尔型、文字型等。

#### 4. 什么是运算符？

运算符如+、-、\*、/等等。

#### 5. 什么是界符？

界符如逗号、分号、括号等等。

## 1.5 实验总体步骤

### 1.5.1 创建编码表

词法分析器的输出是单词序列，在 5 种单词种类中，关键字、运算符、分界符都是程序设计语言预先定义的，其数量是固定的。而标识符、常数则是由程序设计人员根据具体的需要按照程序设计语言的规定自行定义的，其数量可以是无穷多个。编译程序为了处理方便，通常需要按照一定的方式对单词进行分类和编码，在此基础上，将单词表示成二元组的形式（类别编码，单词值）。

单词名称	类别编码	单词值
int	1	-
return	2	-
=	3	-
,	4	-
Semicolon	5	-
+	6	-
-	7	-
*	8	-
/	9	-
(	10	-
)	11	-
id	51	内部字符串
IntConst	52	整数值
.....	.....	.....
布尔常数	80	0 或 1
字符串常数	81	内部字符串

### 1.5.2 创建类 C 语言文法

多数程序语言单词的词法都能用正则文法来描述，基于单词的这种形式化描

述会给词法分析器的设计与实现带来很大的方便，支持词法分析器的自动构造，比如 Flex、ANTLR 词法分析器生成工具。

### 1. 正则文法表示

$G=(V,T,P,S)$ ，其中  $V=\{S,A,B,Cdigit,no\_0\_digit,char\}$ ,  $T=\{\text{任意符号}\}$ ,  $P$  定义如下

约定：用 digit 表示数字：0,1,2,...,9; no\_0\_digit 表示数字：1,2,...,9;

用 letter 表示字母：A,B,...,Z,a,b,...,z,\_

标识符： $S \rightarrow \text{letter } A$            $A \rightarrow \text{letter } A | \text{digit } A | \varepsilon$

运算符、分隔符： $S \rightarrow B$      $B \rightarrow = | * | + | - | / | ( | )$ ;

整常数： $S \rightarrow no\_0\_digit B$            $B \rightarrow digit B | \varepsilon$

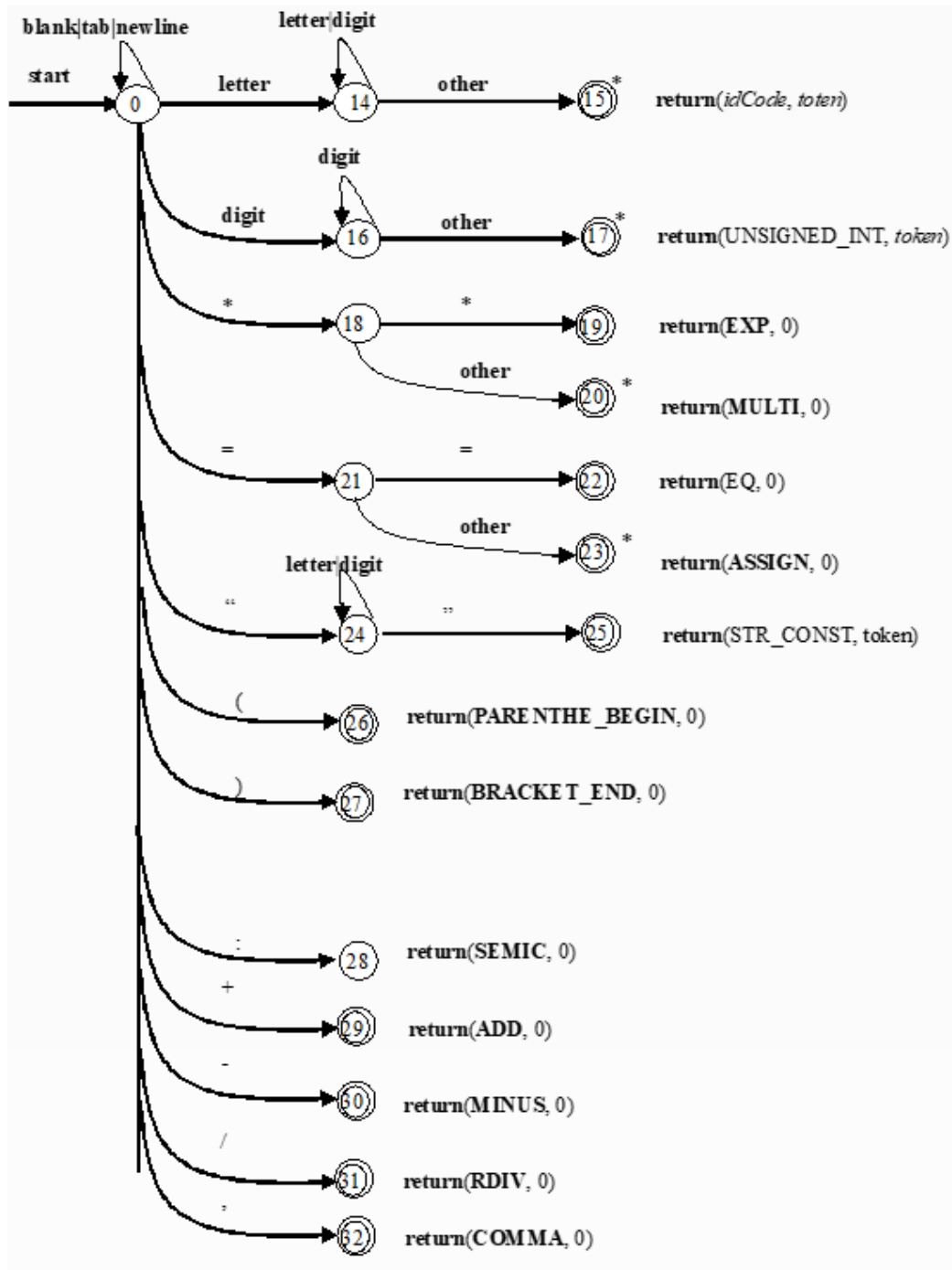
字符串常量： $S \rightarrow "C"$

字符常量： $S \rightarrow 'D'$

### 1.5.3 有限自动机

有限自动机识别的语言称为正则语言，有限自动机分为确定的有限自动机（DFA）和非确定的有限自动机（NFA）两种。DFA 和 NFA 都可以描述正则语言，DFA 规定只能有一个开始符号，且转移标记不能为空，代码实现较为方便，所以词法分析器使用 DFA 描述词法记号。

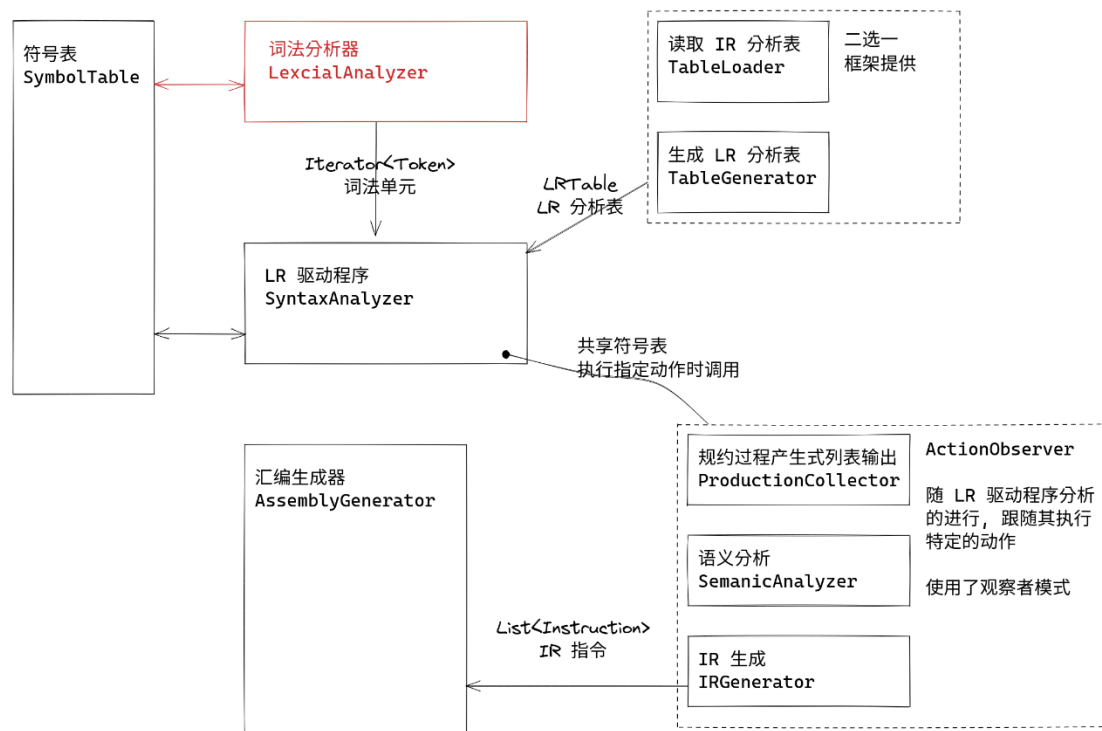
识别各类单词的状态转化图合并图：





## 1.6 实验框架概述

### 1.6.1 编译程序整体框架图



实验一词法分析程序核心类是 `LexcialAnalyzer` 类，该类完成的主要功能包括源程序代码读入，词法分析获取 token 列表以及 token 列表写入文件。

在本次实验中，你将要使用自动机的写法实现一个简单的，特定的词法分析器。具体要实现的词法规则如下：

类别	正则表达式
int	int
return	return
=	=
,	,
Semicolon	;
+	+
-	-
*	*

/	/
(	(
)	)
Id	[a-zA-Z_][a-zA-Z]*
IntConst	[0-9]+

每个词法单元之间可能有任意的空白字符 (\s, 或 [\t\r\n]), 这些字符应当忽略。

你需要根据读入的源语言的代码文本, 生成词法单元迭代器并正确地将源语言中的每个标识符插入到符号表中, 不要求记录每个词法单元的行号, 起始列号, 结束列号, 不要求处理注释。

框架中已包含词法单元 (Token) 以及符号表 (SymbolTable) 的实现。从文件中读取内容、将生成的词法单元输出到文件中、将符号表输出到文件中的过程亦已实现。

## 1.6.2 词法分析程序的输入

### 1. 词法分析程序读入文件

data/in

```
|—— coding_map.csv      # 码点文件（编码表）
|—— input_code.txt       # 输入的代码
```

#coding\_map.csv 码点文件

也就是编码表。

```
1 int
2 return
3 =
4 ,
5 Semicolon
6 +
7 -
8 *
9 /
10 (
11 )
51 id
52 IntConst
```

### #关键字列表

创建关键字列表：关键字可以写入文件中，读文件读入到存储的数据结构；或者直接在程序写入数据结构存储。

注意：这个数据结构模板代码中没有，同学自己创建。

### #input\_code.txt 编译器要解析的源代码

把源程序代码读入到输入缓冲区，根据有限自动机的状态转移识别一个个单词 Token。

```
int result;
int a;
int b;
int c;
a = 8;
b = 5;
c = 3 - a;
result = a * b - (3 + b) * (c - a);
return result;
```

## 1.6.3 词法分析程序的输出

### 1. 词法分析程序输出结果文件

data/out

```
|—— old_symbol_table.txt    # 符号表
|—— token.txt                # 词法单元列表
```

输出的文件内容具体是什么可以参考 `data/std` 目录下的同名文件。

#### # old\_symbol\_table.txt 符号表

编译器用符号表来记录、收集和查找出现在源程序中的各种名字及其语义信息。每当源程序中出现一个新的名字时，则向符号表中填入一个新的表项来记录该名字；当在源程序中发现某个名字的新属性时，要找到该名字所对应的符号表表项，在表中记录其新发现的属性信息。因此，从数据结构的角度来看，符号表是一个以名字为关键字来记录其信息的数据结构，支持插入和查找两个基本操作。在词法分析阶段主要是往符号表里插入名字信息。

符号表的组织结构可以是线性表或者散列表。首先设计者可能会想到选择数组作为符号表的具体物理实现结构，在符号表的规模较小的情况下，线性结构可以满足。当设计的编译器规模较大时，可以考虑引入散列表来提高查找插入的效率。

```
(a, null)
(b, null)
(c, null)
(result, null)
```

#### #token.txt 词法单元列表

词法分析程序通过对输入源程序字符串分割识别，生成一系列独立的单词串，这些单词串包括标识符、关键字、常数、运算符、分界符。

其中关键字、运算符、分界符一符一码，标识符统一一个编码，常量一类一码。词法分析程序识别出单词串以二元组的形式存储，每个二元组由种别码和属性值组成，对于一符一码的单词我们仅记录其种别码，对于标识符、常量这类单词需要同时记录其属性值。

二元组：(种别码, 属性值)

每一个二元组在框架代码中是一个 `Token` 类实例，请参考框架代码。

以下是本实验中所给出的示例输入源程序经过词法分析程序处理后输出的二元组序列，注意(\$,)是在二元组序列末尾加上代表 EOF 的 token。

<i>(int,)</i>	<i>(id,result)</i>
<i>(id,result)</i>	<i>(=,)</i>
<i>(Semicolon,)</i>	<i>(id,a)</i>
<i>(int,)</i>	<i>(*,)</i>
<i>(id,a)</i>	<i>(id,b)</i>
<i>(Semicolon,)</i>	<i>(-,)</i>
<i>(int,)</i>	<i>((,)</i>
<i>(id,b)</i>	<i>(IntConst,3)</i>
<i>(Semicolon,)</i>	<i>(+,)</i>
<i>(int,)</i>	<i>(id,b)</i>
<i>(id,c)</i>	<i>(,)</i>
<i>(Semicolon,)</i>	<i>(*,)</i>
<i>(id,a)</i>	<i>((,)</i>
<i>(=,)</i>	<i>(id,c)</i>
<i>(IntConst,8)</i>	<i>(-,)</i>
<i>(Semicolon,)</i>	<i>(id,a)</i>
<i>(id,b)</i>	<i>(,)</i>
<i>(=,)</i>	<i>(Semicolon,)</i>
<i>(IntConst,5)</i>	<i>(return,)</i>
<i>(Semicolon,)</i>	<i>(id,result)</i>
<i>(id,c)</i>	<i>(Semicolon,)</i>
<i>(=,)</i>	<i>(\$,)</i>
<i>(IntConst,3)</i>	
<i>(-,)</i>	
<i>(id,a)</i>	
<i>(Semicolon,)</i>	

## 1.7 实验框架设计简介

### 1.7.1 数据结构说明

#### 1. Term

该类为所有文法符号（终止符与非终止符）的基类。

#### 2. Token

词法单元的实现，词法单元（Token）是词法分析的结果。

词法分析从源程序文件的文本流中识别出字符，将一个或多个字符根据特定的规则识别出来单独的单词，由这些单词本身和单词的种别码组成词法单元。

一个词法单元逻辑上由两个部分组成：词法单元的类型（种别码），词法单元的单词文本。

```
class Token {  
    private TokenKind kind; //种别码  
    private String text; //单词值  
}
```

然而, 对于大部分简单的词法单元而言, 并不需要保存词法单元代表的文本。在这种情况下, 我们将此 `Token` 的 `text` 设置为"" (空字符串), 而不是 `null`。为了方便构造这两种 `Token`, 我们提供了两个静态方法以供使用:

```
class Token {  
    public static Token simple(TokenKind kind) {  
        return new Token(kind, ""); //不需要保存单词值  
    }  
  
    public static Token normal(TokenKind kind, String text) {  
        return new Token(kind, text); //需要保存单词值  
    }  
}
```

在输出二元组序列末尾, 添加代表 EOF 的二元组单元(`$,` ):

```
    public static Token eof() {  
        return new Token(TokenKind.eof(), "");  
    }
```

### 3. TokenKind

词法单元类别（编码表）。

词法单元的类别一般实现为一枚举, 将所有可能的词法单元类型硬编码到代码中。你也可以在运行时从一个文件中读入所有可能的词法单元类别。

`coding_map.txt` 文件:

```
1 int
2 return
3 =
4 ,
5 Semicolon
6 +
7 -
8 *
9 /
10 (
11 )
51 id
52 IntConst
```

考虑到代码可读性, 我们最终选择使用某个词法单元类别的名字来唯一标识该类别, 而将所谓“码点”作为词法单元类别的附加信息, 仅供参考使用。于是我们的词法单元类别设计如下:

```
class TokenKind{
    private String id;
    private int code;
}
```

鉴于使用字符串作为唯一标识符有极大的因为拼写错误而出错的可能, 我们会在编译器运行时从 `coding_map.txt` 中读入所有可能的词法单元类别标识符, 然后在每次构造词法单元类别时进行检查:

```

class TokenKind {
    public static TokenKind fromString(String id) {
        // 检查 id 是否被允许使用 (即是否在 coding_map.txt 中有定义)
        if (allowed == null || !allowed.containsKey(id)) {
            throw new RuntimeException("Illegal Identifier");
        }

        return allowed.get(id);
    }

    // 禁止外界构造新的 TokenKind
    private TokenKind(String id, int code) {
        // ...
    }
}

```

同时, 为了进一步方便 Token 的构造, 我们也为 simple 和 normal 提供了直接使用词法单元类别标识符的重载版本:

```

class Token {
    public static Token simple(String tokenKindId) {
        return simple(TokenKind.fromString(tokenKindId));
    }

    public static Token normal(String tokenKindId, String text) {
        return normal(TokenKind.fromString(tokenKindId), text);
    }
}

```

另外, TokenKind 作为 Term 的子类, 代表文法中的非终结符。详见实验二指导书。

#### 4. SymbolTable

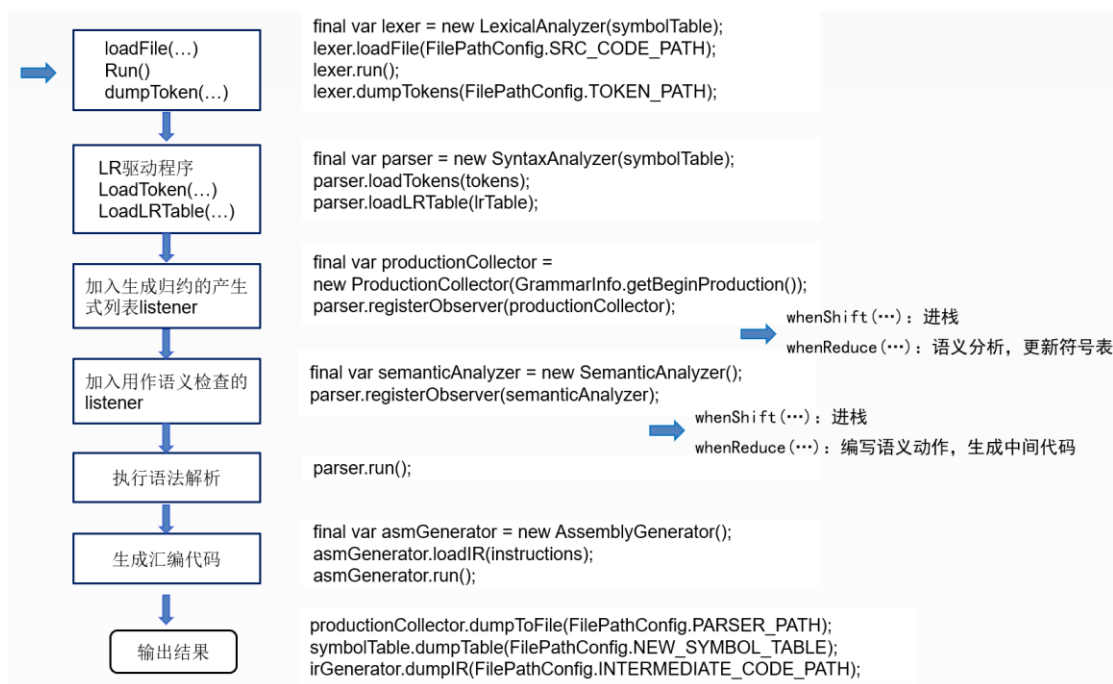
##### 符号表

在实验一中, 你仅需要在读到一个标识符时 (即生成类别为 id 的词法单元时), 检测符号表中是否已含有该标识符, 若无则向符号表添加该标识符即可。注意: 该函数在 SymbolTable 并未定义, 需要自行编写。



```
// 在你的实现的某个地方...
if (!symbolTable.has(identifierText)) {
    symbolTable.add(identifierText);
}
```

## 1.7.2 主程序流程说明



以上是模版代码项目中 Main.java 文件中 main 函数代码流程, 具体代码请大家参考. \template\src\cn\edu\hitsz\compiler\Main.java 文件。

main 函数流程说明如下:

- ①构建符号表, 用来存储源语言中的标识符;
- ②开始词法分析, loadFile 方法用来 load 源语言源代码; 接着执行 run 方法, 也就是词法分析过程, dumpTokens 方法把词法分析的结果 tokens 输出到文件, dumpTable 方法把存储到符号表中的标识符输出到文件。
- ③读取 LR 分析表 (LR 分析表可以用编译工作台工具生成), 框架代码中也提供了 LR 分析表生成的代码, 如需使用框架代码生成请自行启用相关代码。
- ④加载 LR 分析驱动程序; loadTokens 方法读入词法分析的结构 tokens, loadLRTable 读入 LR 分析表, 这两个方法的定义在语法分析程序

SyntaxAnalyzer.java 中，需要开发者自己实现。

⑤框架使用了观察者设计模式，语法分析中需要注册 **ProductionCollector** 用于收集语法分析中归约产生式的观察者接口，执行 **run** 方法实现语法分析。注意：在完成实验二时，需要先把加入语义检查和 **IR** 生成的观察者相关代码注释起来，否则因为这其中有一些代码暂未实现，会导致报异常。

```
// 加入用作语义检查的 Observer
final var semanticAnalyzer = new SemanticAnalyzer();
parser.registerObserver(semanticAnalyzer);

// 加入用作 IR 生成的 Observer
final var irGenerator = new IRGenerator();
parser.registerObserver(irGenerator);
```

⑥框架使用了语法制导翻译的语义分析和中间代码生成，所以语义分析应该是和语法分析过程同步，语法分析每执行一次归约将该次归约的产生式完成语义分析。同样使用观察者模式，注册语义检查和 **IR** 生成的观察者。

⑦**parser.run** 方法执行，完成语法分析和语义分析，在对应观察者中收集归约产生式和完成语义检查以及 **IR** 生成。

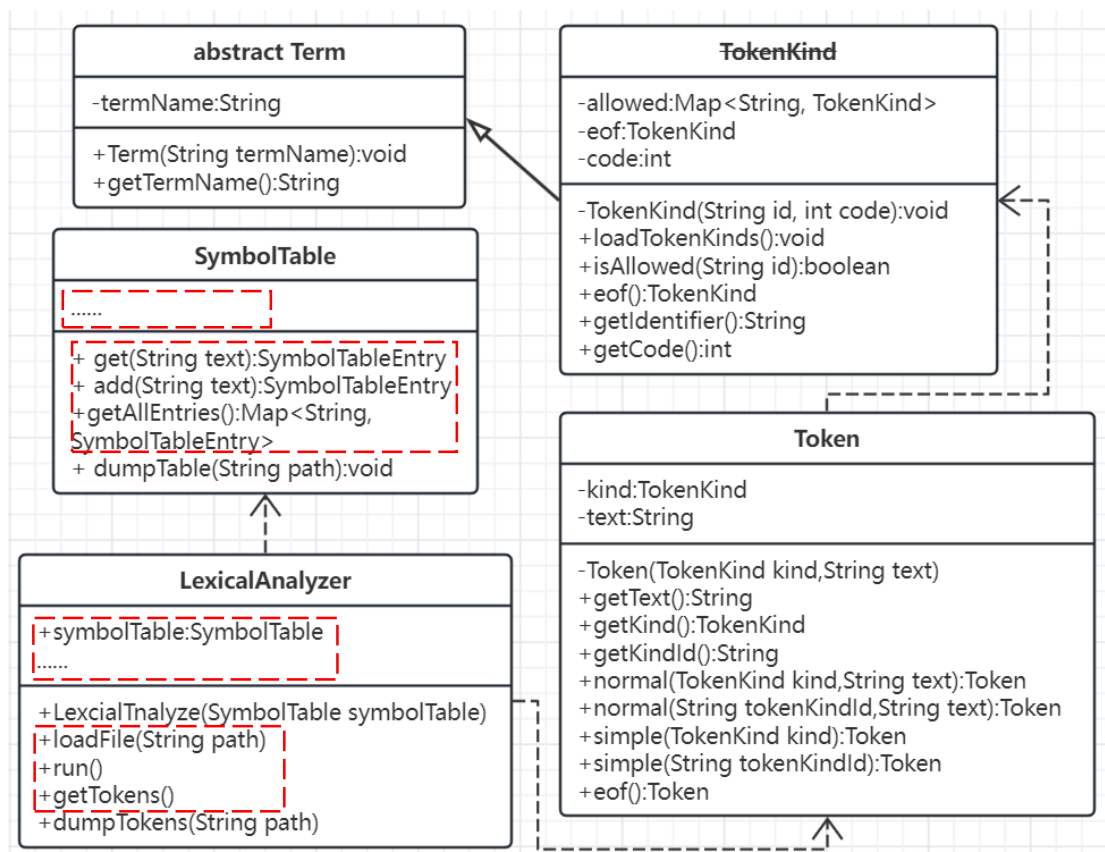
⑧在各观察者中输出结果；

```
//输出语法分析过程中归约的产生式列表
productionCollector.dumpToFile(FilePathConfig.PARSER_PATH);
//输出语义分析语义检查后更新的符号表
symbolTable.dumpTable(FilePathConfig.NEW_SYMBOL_TABLE);
//输出中间代码
irGenerator.dumpIR(FilePathConfig.INTERMEDIATE_CODE_PATH);
```

⑨**loadIR** 把生成的中间代码 **load** 到程序中，执行 **run** 方法，生成汇编代码。

## 1.8 实验框架代码说明

### 1.8.1 词法分析程序框架代码 UML 图



以上是词法分析程序相关类 UML 图，相关说明：

①LexicalAnalyzer 类：这个类是词法分析核心类，开发者需要完成的有 loadFile, run、getTokens 三个函数。loadFile 方法把文件中的源语言源代码 load 到程序中，可以直接采用完整读入的方法读入存储为一个字符串。run 方法执行此番分析，返回一个 token 序列，开发者需要自定义一个列表来存储 token 序列。getTokens 方法从词法分析过程中获取 Token 列表。

②SymbolTable 类：符号表类；词法分析中，当识别到一个标识符，需要把这个标识符名字添加到符号表中。在 SymbolTable 类中，开发者需要实现 has、add、get、getAllEntries 这四个方法，在符号表类中每一个标识符占一行（包括标识符名字和其他属性）也就是一个符号表条目。开发者需要自定义用来存储这些条目的数据结构，可以是一个集合，集合的数据类型应该包含标识符名字和其他属性，其中标识符属性可以使用 SymbolTableEntry 这个模版代码设计好的数据结构。

③Token 类：每一个 token 对象就是一个单词的包装，词法分析结果是一个 token 序列，所以在①中存储词法分析输出结果的数据结构可以使用 Token 类。每一个 token 对象包含单词值本身也就是 text 文本和种别码，getText 方法返回单词本身，getKind 返回 token 的类型，getKindId 返回 token 类型的文本表示；

## 1.8.2 主要类说明

### 1. LexicalAnalyzer 类

```
/**
 * 从给予的路径中读取并加载文件内容
 * @param path 路径
 */
public void loadFile(String path){
    // TODO: 词法分析前的缓冲区实现
    //可自由实现各类缓冲区
    // 或直接采用完整读入方法
    throw new NotImplementedException();
}
```

```
/**
 * 执行词法分析, 准备好用于返回的 token 列表 <br>
 * 需要维护实验一所需的符号表条目, 而得在语法分析中才能确定的符号表条目的成员可以
    先设置为 null
 */
public void run(){
    // TODO: 自动机实现的词法分析过程
    throw new NotImplementedException();
}
```

```
/**
 * 获得词法分析的结果, 保证在调用了 run 方法之后调用
 *
 * @return Token 列表
 */
public Iterable<Token> getTokens() {
    // TODO: 从词法分析过程中获取 Token 列表
    // 词法分析过程可以使用 Stream 或 Iterator 实现按需分析
    // 亦可以直接分析整个文件
    // 总之实现过程能转化为一列表即可
    throw new NotImplementedException();
}
```

## 2. SymbolTable 类:

```
/**
 * 获取符号表中已有的条目
 *
 * @param text 符号的文本表示
 * @return 该符号在符号表中的条目
 * @throws RuntimeException 该符号在表中不存在
 */
public SymbolTableEntry get(String text) {
    throw new NotImplementedException();
}
```

```
/**
 * 在符号表中新增条目
 *
 * @param text 待加入符号表中的新符号的文本表示
 * @return 该符号在符号表中对应的新条目
 * @throws RuntimeException 该符号已在表中存在
 */
public SymbolTableEntry add(String text) {
    throw new NotImplementedException();
}
```

### 1.8.3 补充说明

框架代码下载地址

<https://gitee.com/hitsz-cslab/Compiler/releases/tag/2022F.0.1>

在线指导书

<https://hitsz-cslab.gitee.io/compiler/>

说明：

1. 框架模板仅支持声明语句、赋值语句和简单算术表达式。同学想开发的编译器支持更多语法功能需自行修改以下文件：
  - ✧ 源代码文件：input\_code.txt;
  - ✧ 编码表文件：coding\_map.txt;
2. 注意事项：
  - ✧ 禁止使用 java.util 提供的正则表达式的文本处理工具来完成词法分析；
  - ✧ 禁止直接查找空格来分割单词；