

# CS 131 Compilers: Discussion 1: Lexical Analysis

杨易为 季杨彪 尤存翰

{yangyw,jiyb,youch}@shanghaitech.edu.cn

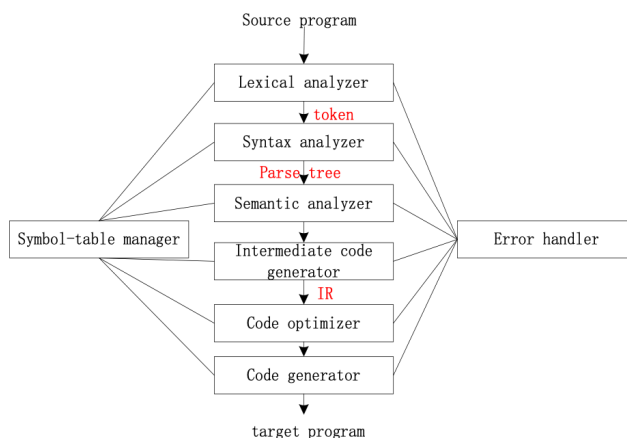
2021 年 2 月 12 日

## 目录

## 1 Welcome to the Compiler's World

### 1.1 Some General Views

We've covered on class that the compiler includes lexical analysis, in which we get tokens. This process is called scanner of the program.



To make detail definition of compiling a program into a machine recognizable text and discuss whether kind of language is "readable" by a Turing Machine or so-called computer, we introduce regular expressions. Formally, we have kleene closure and possitive closure of a language to describe a language. Then NFA, DFA, DFA minimization are tools to accept the language. Regex have some limitation of describing a language and how to evaluate the power of it? Context Free Grammer is introduced and its solution LL Parsing and LR Parsing.

Till now, we have made the program into a Intermediate Representations, which in our project is Abstract Syntax Tree. We'll then make some type checking on it. This takes into the definition of Operational Semmantics. For modern language like LLVM or those who compiled into LLVM IR, they have a Intermediate Representations where passes can be cast on to optimize the code. Those optimizations are basically machine unrelated. Other intermediate representations include 3 address code and Single Static Code. Some of them are more low level and others are very similar to the AST, which depends how heavy the analysis and the optimization can cast on. We'll cover intra-basic-block optimization, dataflow analysis, security problems like Control Hijacking Attack. Language may alloacate resources on stack or heap, we can do runtime garbage collections to make the resources reallocatable.

The eventual process of a compiler is Code Generator to generate machine code, we'll cover transform the IR into the riscv backend that we've covered on Computer Architechture classes. Wish you a good luck at Compiler class.

## 1.2 Why Both Parser and Lexer?

Julia is language with no Lexer and some of the parser , it generates Single Parser or Single Lexer is not sounded to compiled a well defined language.

## 1.3 How can a simple regex work in your project?

1. Any Search with Regex option.
2. You can write a plugin for a new language like Julia to support highlighting in vim/vscode/jetbrains.

## 1.4 Stages of Compilations

### 1.4.1 Convert the C++ source code into a sequence of tokens.

Compilers are conceptually broken down into multiple phases, each of which carry out transformations on a source program towards the compiler's ultimate goal; the generation of a target program. The result of each phase is an intermediate representation that facilitates the implementation of subsequent phases (or, if it is the last phase, the result is the target program). The first phase in a compiler, the lexical analyzer or lexer, reads an input stream of characters and converts it into a stream of tokens to be sent to the parser for syntactic analysis.

```

1 // n! = 1 * 2 * 3 * .... * (n-1) * n
2 int factorial(int n) {
3     int result = 1;
4     for(int i = 2; i <= n; i++)
5         result *= i;
6     return result;
7 }
```

**Answer:**

### 1.4.2 The rest of the compilations.

What are the rest of the stages of compilation? List them in order and write a brief description about what they and what they take as input and what they output.

**Answer:**

## 2 Regular Expressions

The lexer needs to scan and identify finite character sequences that match a pattern corresponding to a particular token (the matching strings are also called lexemes). Patterns on finite strings can be codified using regular expressions, each of which specify what is called a regular language. Kleene's theorem states

that any regular language can be recognized by a finite state automaton (FSA) and any language that is recognized by an FSA is regular.

## 2.1 What language is denoted by each of the following regular expressions?

Try writing down a few simple strings and give a concise description of the language.

### 2.1.1 $[a-zA-Z][a-zA-Z0-9]^*$

**Answer:**

### 2.1.2 $((\epsilon \mid a)b^*)^*$

**Answer:**

## 2.2 Write a regular expression for the following languages.

### 2.2.1 Uppercase Letters

All strings of uppercase letters, where the letters are in ascending lexicographic order (empty string allowed).

**Answer:**

### 2.2.2 ShanghaiTech Teacher's email addresses

xxxx@shanghaitech.edu.cn, different error handler for different inputs?

**Answer:**

Suppose you have a large spreadsheet, each row of which has a first name, last name, and email column. You want to find all people who use an email address of the form first name + last name @ shanghaitech.edu.cn. Write a regular expression using backreferences that describes this pattern.

**Answer:**

### 2.2.3 Months and Dates

The following four strings: October 8th, October 8, Oct 8th, Oct 8 (Be as concise as possible with your regex)

**Answer:**

### 2.2.4 Ones and Zeros

Even binary numbers without leading 0.

**Answer:**

