# CS 536 Announcements for Monday, February 5, 2024

**Programming Assignment 1**
- symbol table files due Thursday, Feb. 8 by 11:59 pm

**Homework 0**
- available in schedule
- practice with DFAs, regular expressions

**Homework 1**
- available tomorrow
- practice with NFA→DFA translation, JLex

**Last Time**
- non-deterministic FSMs
- equivalence of NFAs and DFAs
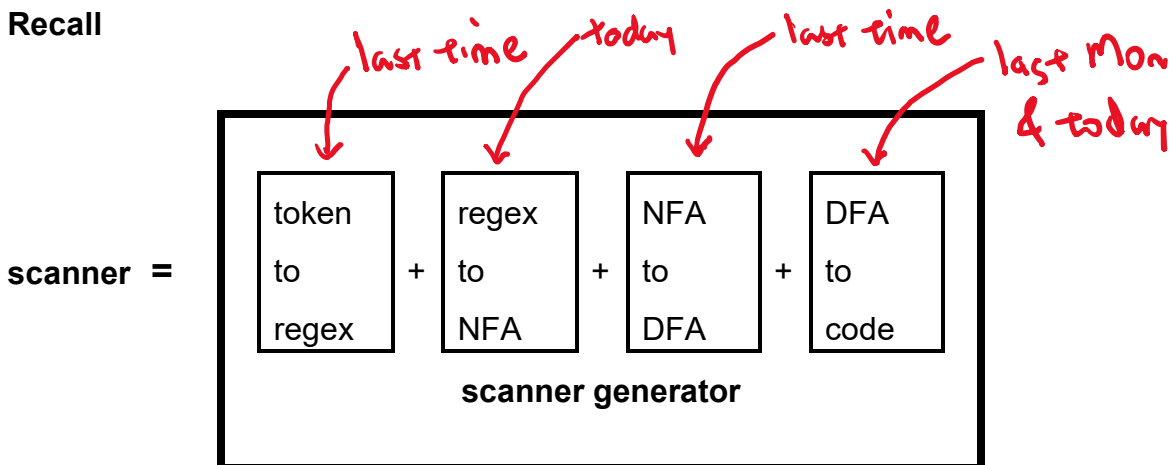- regular languages
- regular expressions

**Today**
- regular expressions → DFAs
- language recognition → tokenizers
- scanner generators
- JLex

**Next Time**
- CFGs

**Recall**

*last time* *today*   *last time*   *last Mon & today*

**scanner =**

| token to regex | + | regex to NFA | + | NFA to DFA | + | DFA to code |
|---|---|---|---|---|---|---|

**scanner generator**

# From regular expressions to NFAs

**Overview of the process**

- Conversion of literals and epsilon $\longrightarrow$ simple FAs
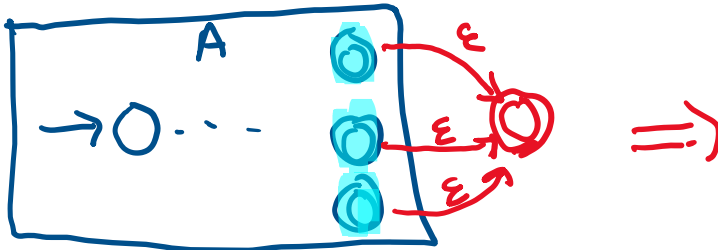- Conversion of operators
  - Convert operands to NFAs
  - join NFAs

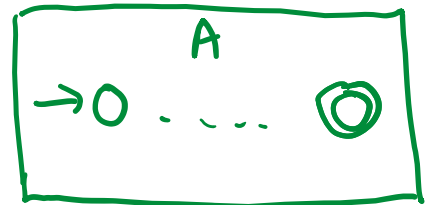## Regex to NFA rules

**Rules for operands**

literal 'a'   $\rightarrow \bigcirc \xrightarrow{a} \circledcirc$

epsilon $\varepsilon$   $\rightarrow \bigcirc \xrightarrow{\varepsilon} \circledcirc$
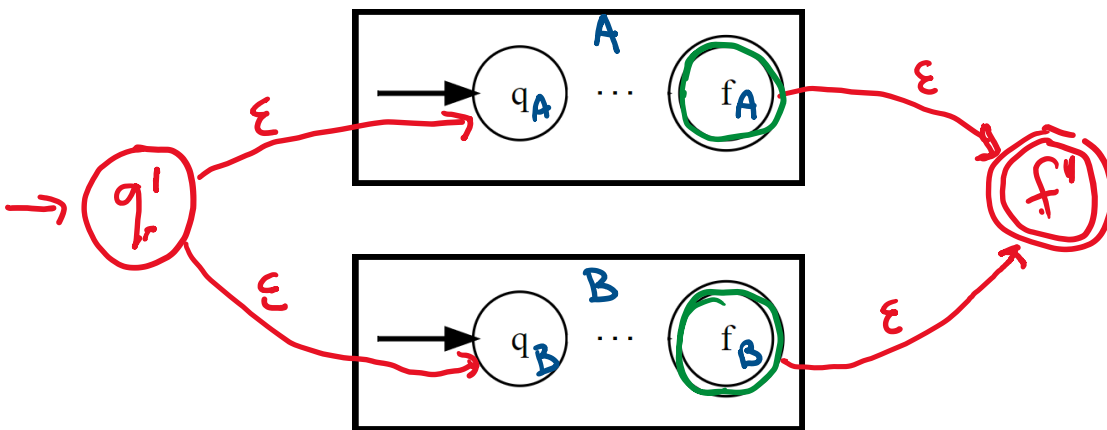
**Suppose A is a regex with NFA:**

Convert so only 1 final state

make these non-final

**Rules for alternation  A|B**

make $f_A, f_B$ non-final

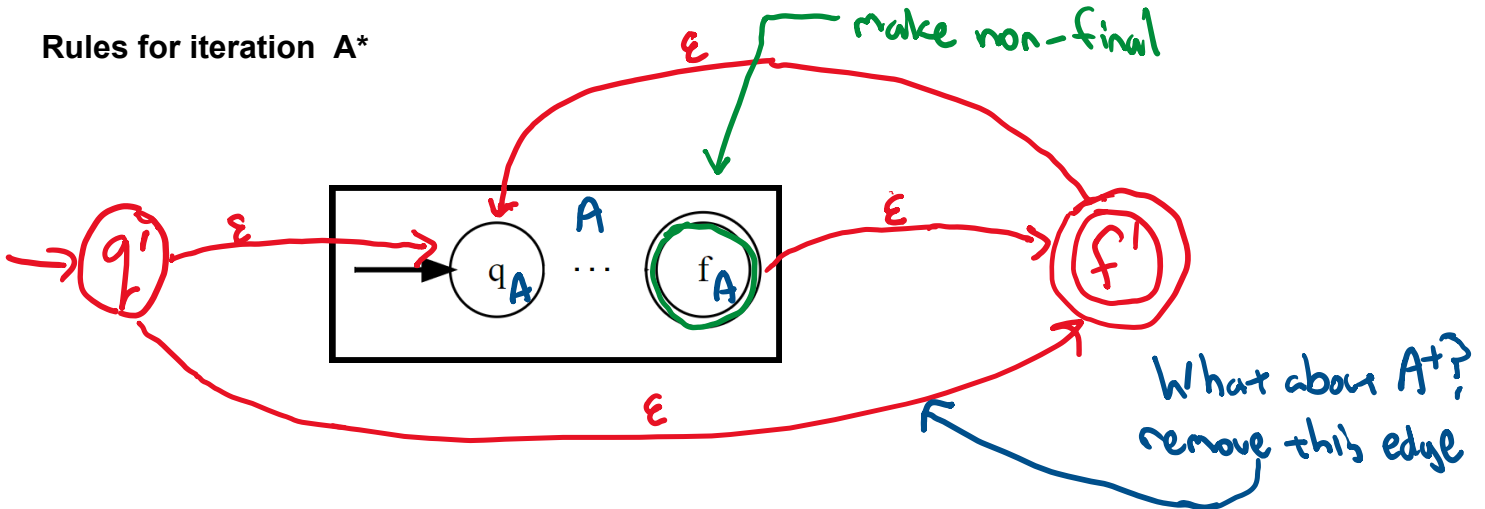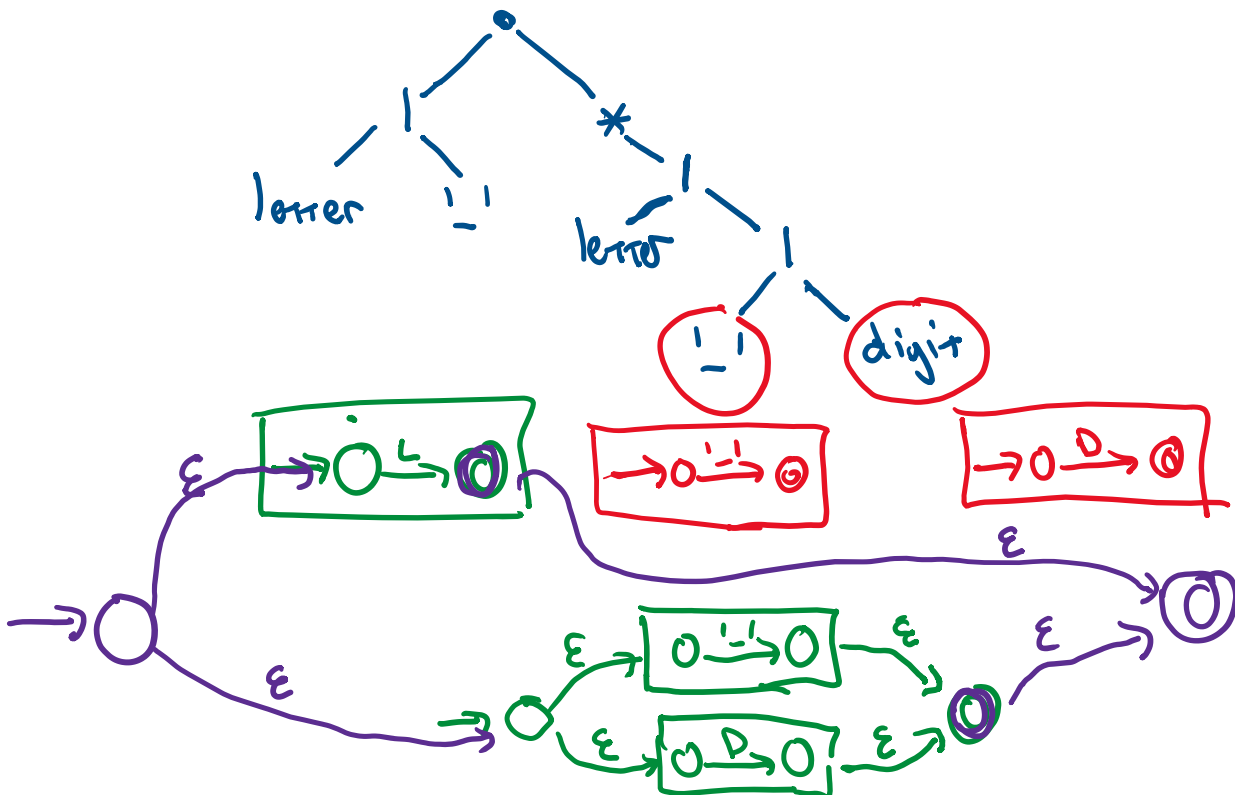# Regex to NFA rules

**Rules for catenation  A.B**

make non-final



**Rules for iteration  A\***

make non-final

What about A+?
remove this edge



# Tree representation of a regex

**Consider regex:**  ( letter | '_' ) ( letter | '_' | digit )*

# Regex to DFA

We now can do:

$$regex \rightarrow NFA_{w/\varepsilon} \rightarrow NFA_{w/o\varepsilon} \rightarrow DFA$$

We can add one more step: **optimize DFA**

**Theorem:** For every DFA *M*, there exists a unique equivalent smallest DFA *M\** that
recognizes the same language as *M*.

↳ fewest # of states

**To optimize:**

- remove unreachable states  → can't get to from start state
- remove dead states  → can't get to a final state from it
- merge equivalent states
  ↳ same transitions (out) w/ same labels

# But what's so great about DFAs?

**Recall:** state-transition function (δ) can be expressed as a table

➔ very efficient array representation

|       | a     | b     | c     |
|-------|-------|-------|-------|
| $S_1$ | $S_2$ | $S_2$ |       |
| $S_2$ | $S_1$ |       | $S_2$ |

➔ efficient algorithm for running (any) DFA

```
s = start state
while (more input){
    c = read next char
    s = table[s][c]
}
if s is final, accept
else reject
```

# What else do we need?

**FSMs** – only check for language membership of a string

**scanner** needs to
- recognize a stream of many different tokens using the longest match
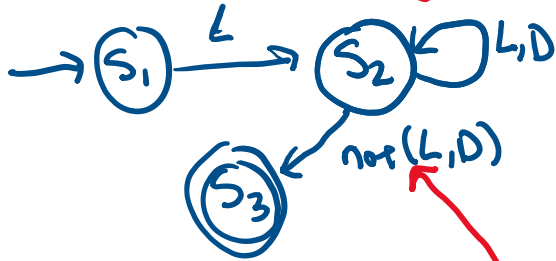- know what was matched

# Table-driven DFA → tokenizer

**Idea:** augment states with actions that will be executed when state is reached

Consider: ( letter )( letter | digit )*



| State | action |
|---|---|
| $S_2$ | return ID ← token |

Problem: **Don't get longest match**

$\boxed{avg = 7}$

| State | action |
|---|---|
| $S_3$ | ~~return ID~~ |

$\boxed{avg = 7}$

To fix:
put back
1 char,
return ID

Problem: **maybe we need this char**

**Actions needed:**
- return a token
- put back a character
- report an error

Also add EOF token,
EOF symbol to alphabet $\Sigma$

# Scanner Generator Example

**Language description:**
consider a language consisting of two statements

- assignment statements: `ID = expr`
- increment statements: `ID += expr`

where `expr` is of the form:

- `ID + ID`
- `ID ^ ID`
- `ID < ID`
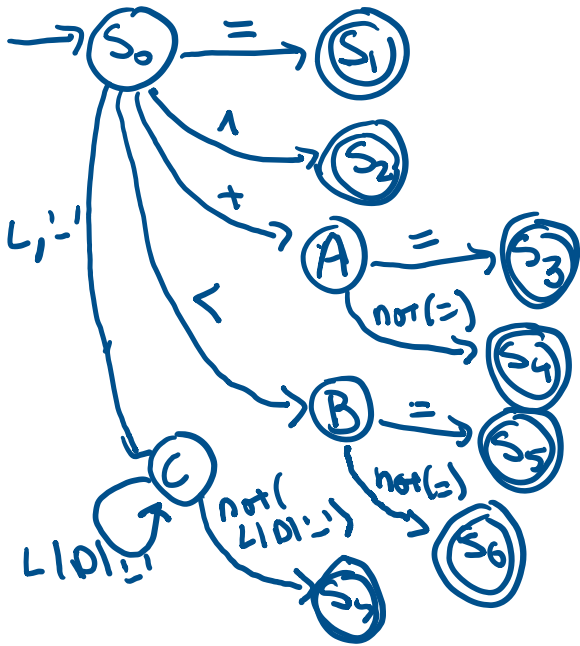- `ID <= ID`

and `ID` are identifiers following C/C++ rules (can contain only letters, digits, and underscores; can't start with a digit)

**Tokens:**

| Token | Regular expression |
|---|---|
| ASSIGN | " = " |
| INCR | " += " |
| PLUS | " + " |
| EXP | " ^ " |
| LESSTHAN | " < " |
| LEQ | " <= " |
| ID | |

$(\text{letter} \mid {}'\_{}')(\text{letter} \mid {}'\_{}' \mid \text{digit})^*$

## Combined DFA



Actions
S₁: return ASSIGN
S₂: return EXP
S₃: return INC
S₄: put 1 back, return PLUS
S₅: return LEQ
S₆: put 1 back, return LESSTHAN
S₇: put 1 back, return ID

## State-transition table

|   | = | + | ^ | < | _ | letter | digit | EOF | none of these |
|---|---|---|---|---|---|--------|-------|-----|---------------|
| S₀ | ret ASSIGN | A | ret EXP | B | C | C |  | ret EOF |  |
| A | ret INC | put 1 back, ret PLUS |  |  |  |  |  |  | → |
| B | ret LEQ | put 1 back, ret LESSTHAN |  |  |  |  |  |  | → |
| C | put 1 back, ret ID |  |  | C | C | C | | put 1 back, ret ID | → |

```
do {
    read char
    perform action / update state
    if (action was to return a token)
        start again in start state
} while not(EOF or stuck)
```

# Lexical analyzer generators
## (aka scanner generators)

Formally define transformation from regex to scanner

Tools written to synthesize a lexer automatically

- Lex : UNIX scanner generator, builds scanner in C
- Flex : faster version of Lex
- JLex : Java version of Lex

## JLex
### Declarative specification (non-procedure)
- you don't tell JLex <u>how</u> to scan / how to match tokens
- you tell JLex <u>what</u> you want scanned (tokens) & what to do when a token is matched

**Input:** set of regular expressions + associated actions

JLex specification

**Output:** Java source code for a scanner

— .jlex eg xyz.jlex

↳ xyz.jlex.java ⅂ compile to get Yylex.class

- ctor: takes input stream as arg
- next_token: return next token of input

### Format of JLex specification
3 sections separated by %%
- user code section
- directives
- regular expression rules

## Regular expression rules section
**Format:** `<regex>{code}` where `<regex>` is a regular expression for a single token
- can use macros from Directives section – surround with curly braces `{ }`
- characters represent themselves (except special characters)
- characters inside " " represent themselves (except \" )
- . matches anything

**Regular expression operators:**  |   *   +   ?   ( )

**Character class operators:**     –      ^      \

# JLex example

```
// This file contains a complete JLex specification for a very
// small example.

// User Code section:  For right now, we will not use it.

%%

DIGIT=        [0-9]
LETTER=       [a-zA-Z]
WHITESPACE=   [\040\t\n]

%state SPECIALINTSTATE

%implements java_cup.runtime.Scanner
%function next_token
%type java_cup.runtime.Symbol

%eofval{
System.out.println("All done");
return null;
%eofval}

%line

%%

({LETTER}|"_")({DIGIT}|{LETTER}|"_")* {
                          System.out.println(yyline+1 + ": ID "
                                 + yytext()); }

"="           { System.out.println(yyline+1 + ": ASSIGN"); }
"+"           { System.out.println(yyline+1 + ": PLUS"); }
"^"           { System.out.println(yyline+1 + ": EXP"); }
"<"           { System.out.println(yyline+1 + ": LESSTHAN"); }
"+="          { System.out.println(yyline+1 + ": INCR"); }
"<="          { System.out.println(yyline+1 + ": LEQ"); }
{WHITESPACE}* { }
.             { System.out.println(yyline+1 + ": bad char"); }
```

## Using scanner generated by JLex in a program

```
// inFile is a FileReader initialized to read from the
// file to be scanned
Yylex scanner = new Yylex(inFile);
try {
    scanner.next_token();
} catch (IOException ex) {
    System.err.println(
            "unexpected IOException thrown by the scanner");
    System.exit(-1);
}
```