# CS 344 - Regular expressions

## Announcements

- Ch 2 this week
- Essay due Friday

# Regular Expressions

## Defined as:

- A character

- The empty string, $\varepsilon$

- 2 regular expressions concatenated

- 2 regular expressions separated by an or (written |)

- A regular expression followed by * (Kleene star — 0 or more occurrances)

# Regular Languages

The class of languages described by a regular expression.

Ex: $\underbrace{0^*}\underbrace{10^*} = L$

language of all strings that contain exactly one 1

$\in$ is a member of

$1 \in L$

$01 \in L$      $11 \notin L$

$00010 \in L$      $0 \notin L$

**Ex:** Give the regular expression for
$\{w \mid w$ begins with a 1 and ends with a 0$\}$

$$1(0|1)^*0$$

**Ex:** $\{w \mid w$ starts with 0 and has an odd length$\}$

$$0((0|1)(0|1))^*$$

# Example: Numbers in Pascal

digit → [0-9]

$$\left[\begin{array}{l} \text{digit} \longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\ \text{unsigned\_int} \longrightarrow \text{digit digit}* \end{array}\right.$$

unsigned_number →
unsigned_int ( ε | . unsigned_int )
( ε | ( ( e | E ) ( + | - | ε ) unsigned_int ) )

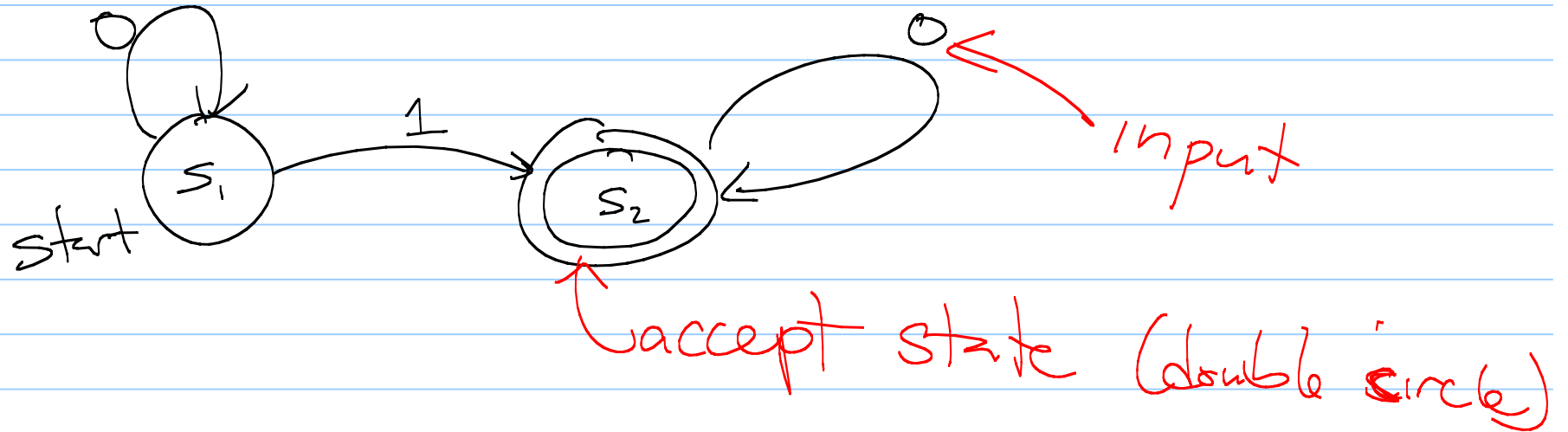# Deterministic Finite Automate (DFA)

Regular languages are precisely the things ~~recognized by DFAs~~.

- A set of states
- input alphabet
- A start state
- A set of accept states
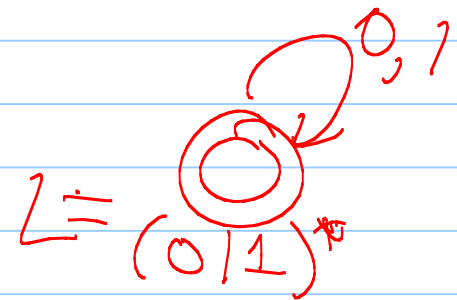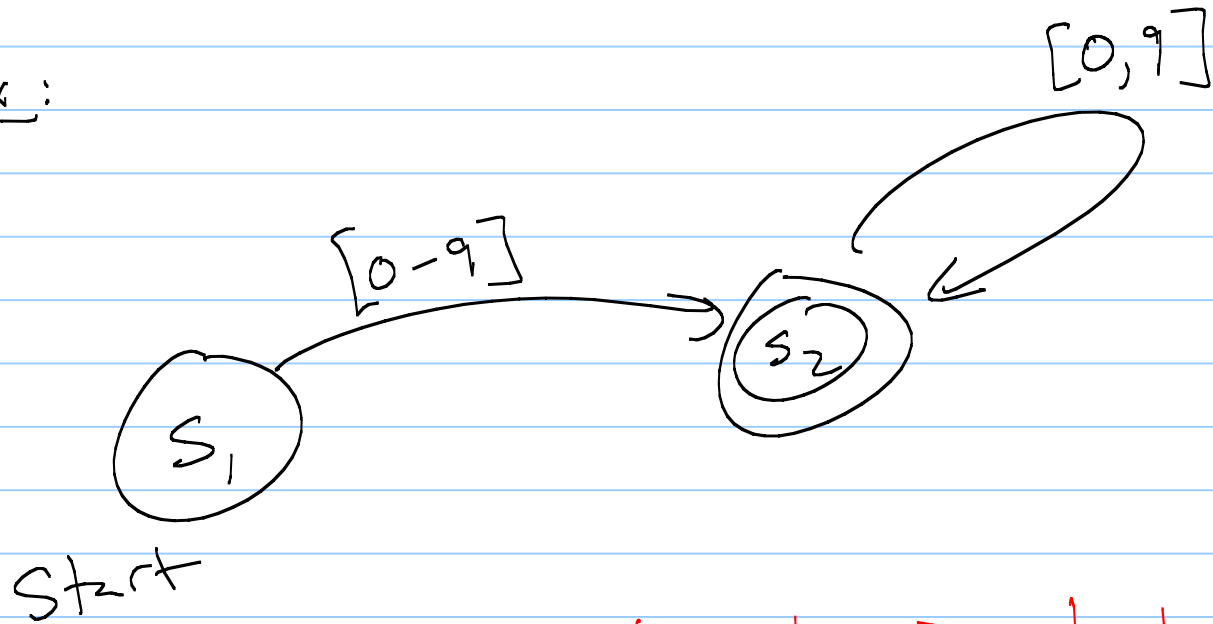- A transition function: given a state & an input, output a new state

# Example:



0010

arrows give transition fn

input

accept state (double circle)

$L$ accepted is $0^* 1 0^*$
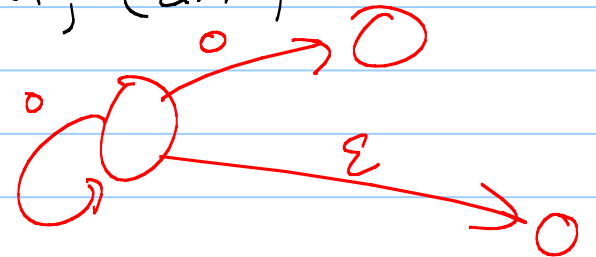
$L = (0/1)^*$

Ex:

[0,9]

[0-9]

S₁    →    S₂

Start

unsigned_int → digit digit*

# Non-deterministic Finite Automata
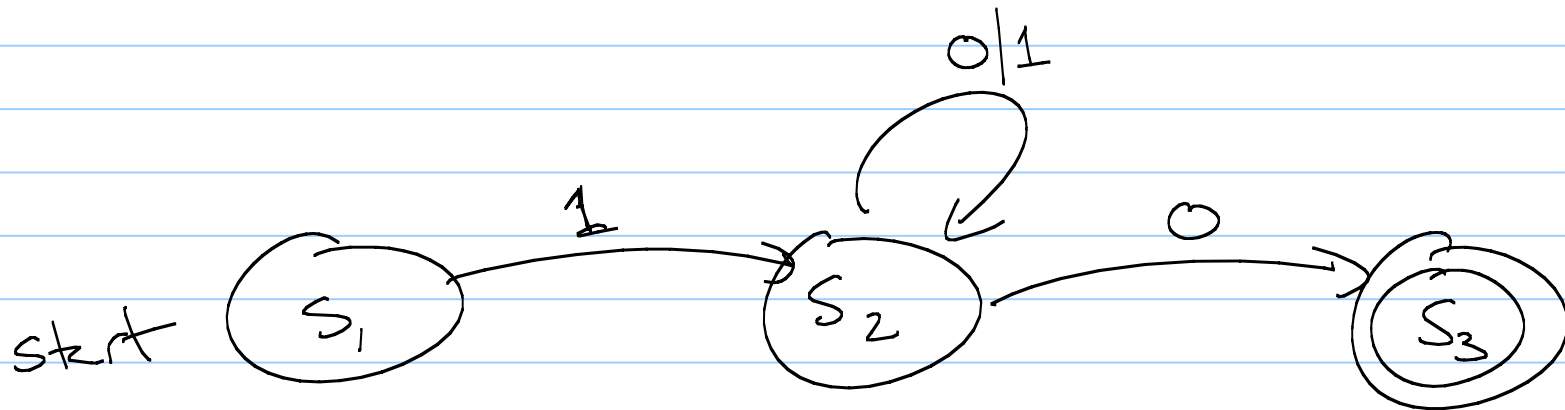
Note: No ambiguity is allowed in DFA's.
So given a state & input, can't
be multiple options.

Also — no $\varepsilon$-transitions.

If we allow several choices to
exist, this is called an NFA.

$$1(0|1)^* 0$$

## Ex :

unsigned_number →
unsigned_int ( ε | . unsigned_int )

Essentially, we can think of an NFA as modeling a parallel set of possibilities (or a tree of them).

Thm: Every NFA has an equivalent DFA.

Both recognize reg. languages.

# Limitations of Regular Expressions

Certain languages are not regular.

Ex: $\{w \mid w$ has an equal number of $0$'s and $1$'s $\}$

Somehow, this needs a type of memory, which regular expressions do not have.

$$0^n 1^n$$

Why do we need this?

Need to "nest" expressions.

Ex: $expr \rightarrow id \mid number \mid -expr \mid$
$(expr) \mid expr \ op \ expr$

$op \rightarrow + \mid - \mid / \mid *$

Regular expressions can't quite do this.

# Context Free Languages - CFLs

Described in terms of productions
   (called Backus-Naur Form, or BNF)

- A set of terminals T

- A set of non-terminals N

- A start symbol S   (a non-terminal)

- A set of productions

Ex: $\{0^n 1^n \mid n > 0\}$

Nonterminals $= \{0, 1\}$

$S \longrightarrow 0 S 1$

$S \longrightarrow \varepsilon$

} productions

Ex: $\{w \mid w$ has an equal number of 0's & 1's$\}$

$$S \longrightarrow 0S1 \qquad ☆$$

$$S \longrightarrow 1S0$$

$$S \longrightarrow \varepsilon$$

Show $001011$ is in this lang:

$$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 001S011$$
use $\varepsilon$

# Expression grammars: Simple calculator

$expr \rightarrow term \mid expr \ add\_op \ term$

$term \rightarrow factor \mid term \ mult\_op \ factor$

$factor \rightarrow id \mid number \mid -factor \mid (expr)$

variable

terminals

$add\_op \rightarrow + \mid -$

$mult\_op \rightarrow * \mid /$

$x + y * z + a$

**Example:** Show how rules can generate $3 + 4 * 5$

expr $\Rightarrow$ expr addop term

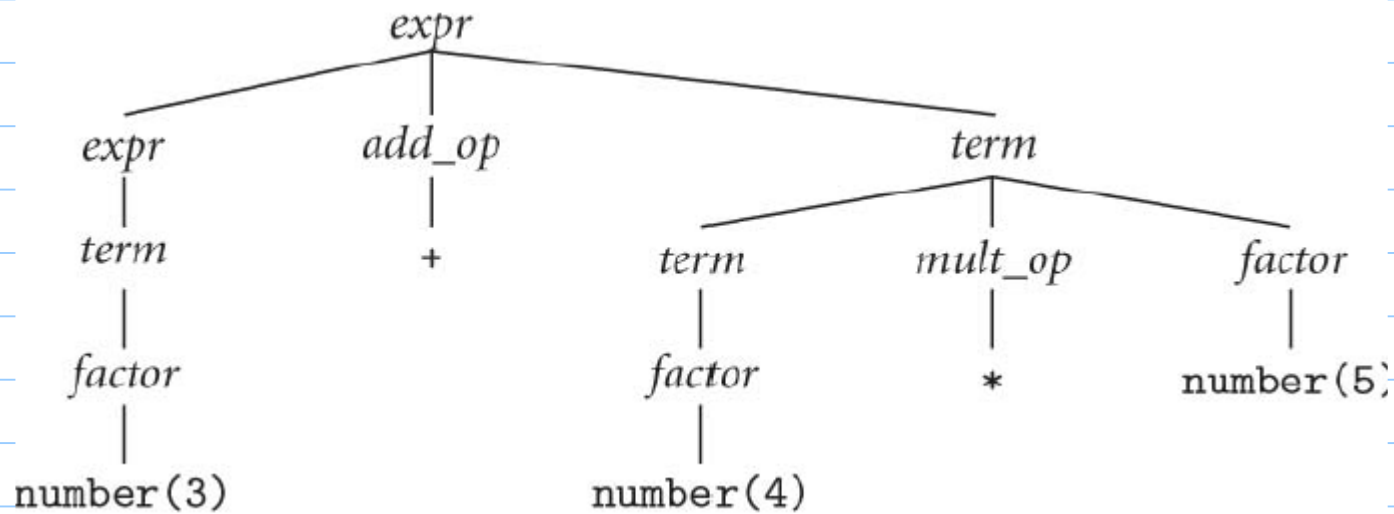{ expr $\rightarrow$ term
term $\rightarrow$ factor
factor $\rightarrow$ num

$\Rightarrow$ number + term

$\Rightarrow$ number + (term * factor)

{ term $\rightarrow$ factor
factor $\rightarrow$ num

{ factor $\rightarrow$ num

num + (num * factor)

# Parse Tree

Ex: 3 + 4 * 5

```
                              expr
                    ┌──────────┼────────────────────────┐
                  expr       add_op                    term
                    │          │              ┌──────────┼──────────┐
                  term         +             term     mult_op     factor
                    │                          │          │          │
                 factor                     factor        *      number(5)
                    │                          │
               number(3)                   number(4)
```

# Another example:
Generate: slope * x + intercept



```
                              expr
                 _____/   |   _____
                /               |               \
            expr             add-op            term
              |                 |                |
            term               +              factor
       ___/   |   \___                          |
      /       |       \                      id = intercept
   term    multop    factor
     |                  |
   factor            id = x
     |
   id = slope
```

**Question:** Can these be ambiguous?

Ex: 10 - 4 - 3

Yes

expr

expr ──── addop ──── term

expr ──── addop ──── term

│ │ │

exp addop term factor

│ │ │ │

⋮ — factor number = 3

│ │

10 number = 4

← not this one "

# Scanning

Recall that the scanner is responsible for

- tokenizing source code

- removing comments

- save text of identifiers, #s, strings

- saving source locations for error messages

# Ex: Calculator Scanner

expr $\rightarrow$ term | expr add_op term
term $\rightarrow$ factor | term mult_op factor
factor $\rightarrow$ id | number | -factor | (expr)
add_op $\rightarrow$ + | -
mult_op $\rightarrow$ * | /

## Add:

id $\rightarrow$ letter (letter | digit)*

comment $\rightarrow$ /* (non-* | * non-/)* */
                    // (non-newline)* newline

assign $\rightarrow$ :=

How to scan / recognize?

# Ad Hoc Approach: Go char by char!

If current $\in \{$ "(", ")", "+", "-", "*" $\}$

    return that symbol

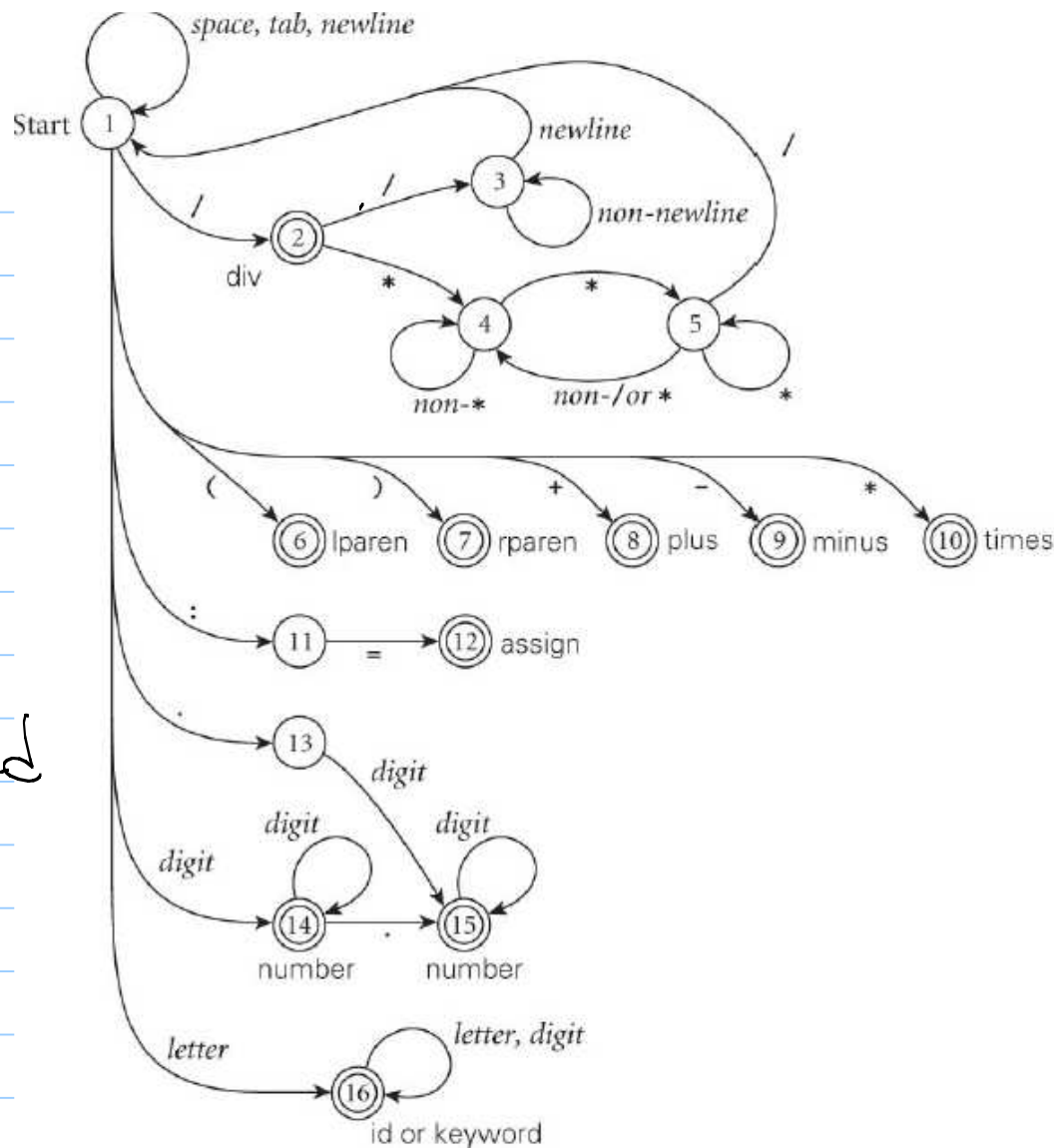If current = ":"

    read next

      if it is =, announce "assign"

        else announce error

if current = "/"

    read next

      if it is "*" or "/"

        read until "*/" or "newline" (resp.)

    else return divide

**Another Way:**

We are essentially running a/ DFA!

Accept states are just places we have reached a token.

## Getting tokens (with DFA)

- Run the machine over & over to get our tokens

Rule: Always take longest possible token

Why?    Ex:    3.14159


Ex: foobar