# CS 344 — Scanning

## Announcements

- Essay due on Friday
- Next HW will be up by Friday

$$d^+ = d(d)^*$$

# Last time: Regular expressions

- 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 occurrences)

+ operator

$(d)^+$

1 or more occurrences

Ex: Give the regular expression for
{w| w begins with a 1 and
ends with a 0}

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

Ex: {w| w starts with 0 and has
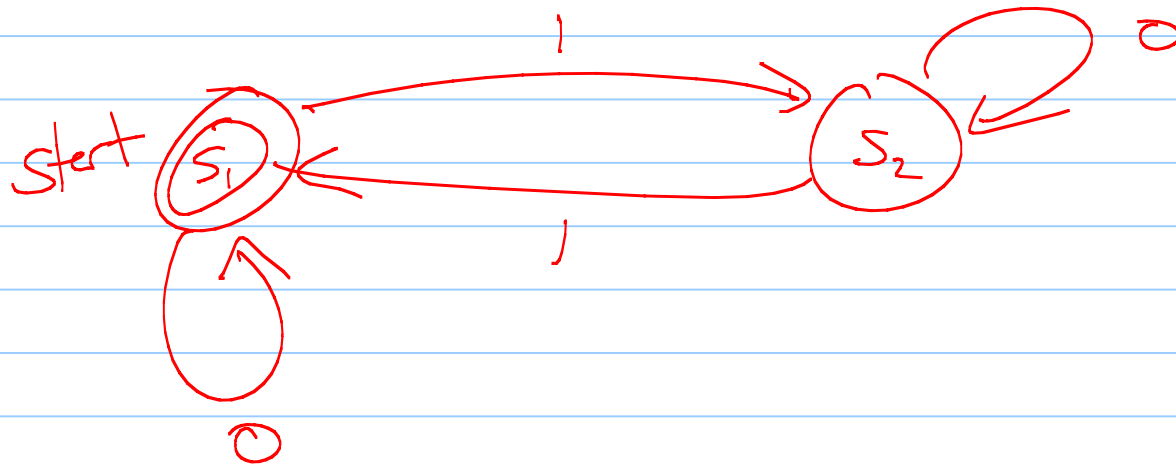an odd length}

$$0((0|1)(0|1))^*$$
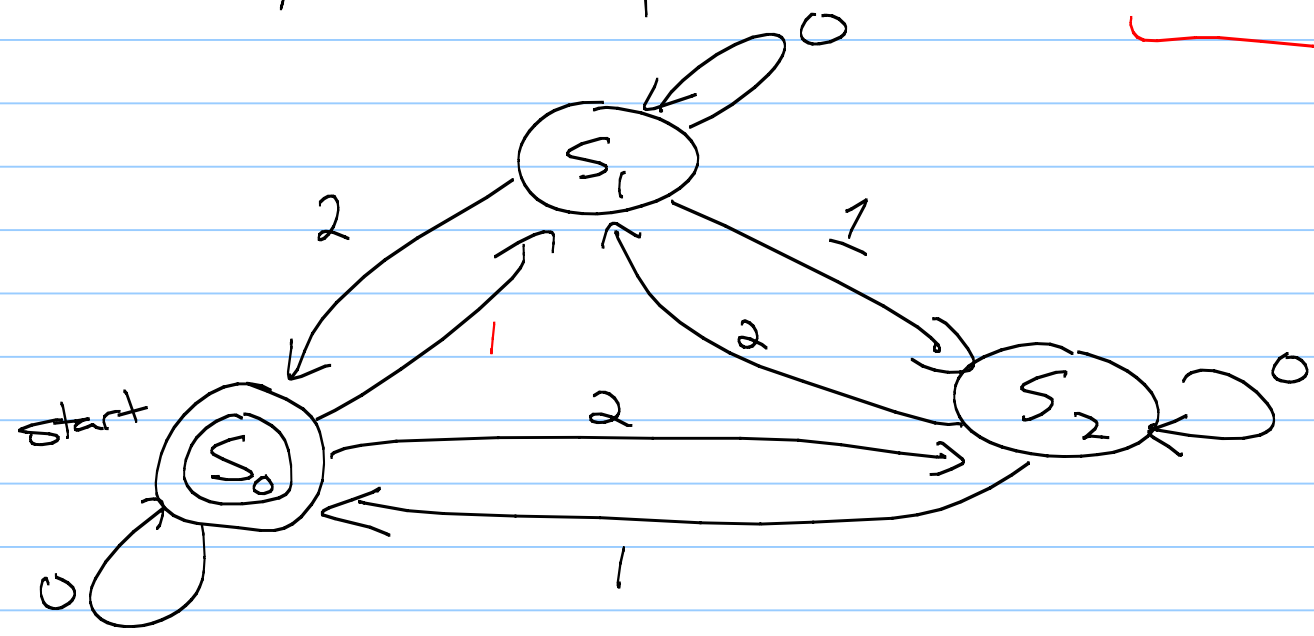
# Deterministic Finite Automate (DFA)

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

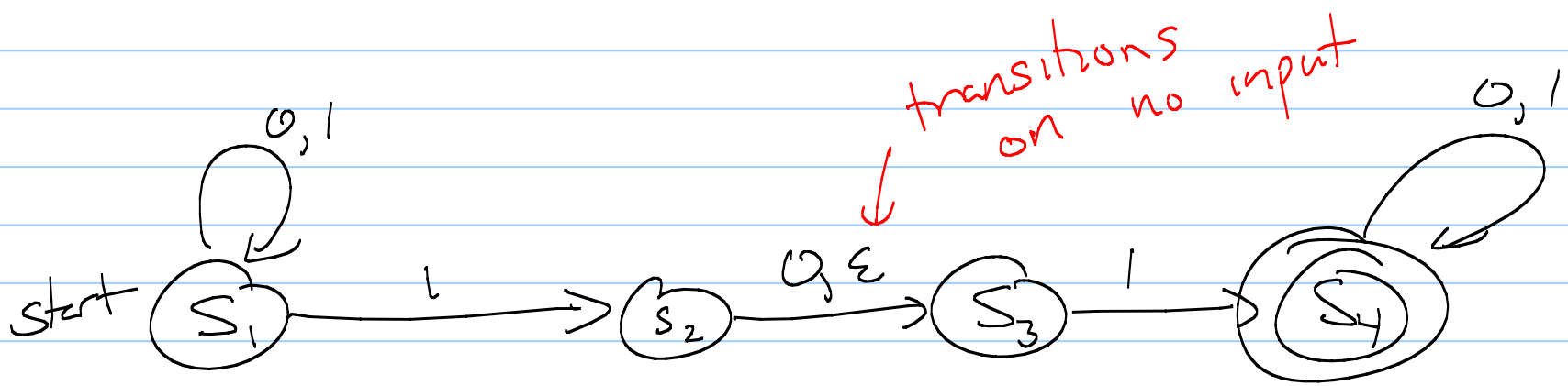# Ex: String of 0's & 1's:
accept if number of 1's is
<u>even</u>

Ex: 3 symbol alphabet: $\{0, 1, 2\}$



computing mod 3
accepting if sum of word is 0 mod 3

# NFAs: DFAs w/ ambiguity

start $\rightarrow$ $S_1$ $\circlearrowleft$ 0,1

$S_1$ $\xrightarrow{1}$ $S_2$ $\xrightarrow{0, \varepsilon}$ $S_3$ $\xrightarrow{1}$ $S_4$ $\circlearrowleft$ 0,1
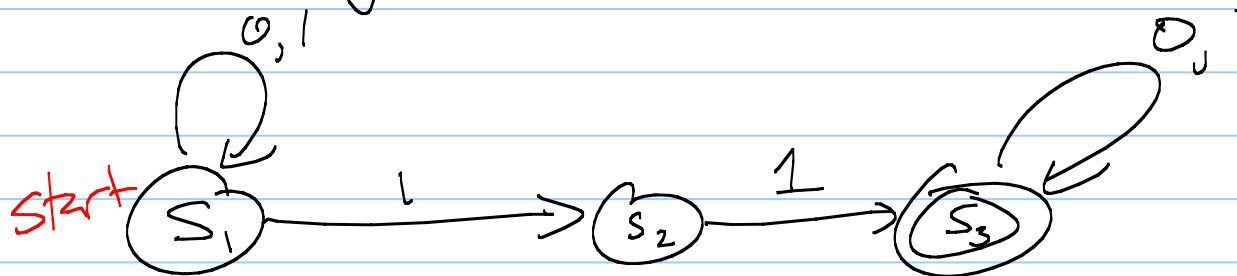
transitions on no input

if read in a 1, could go to $S_1$ or $S_2$

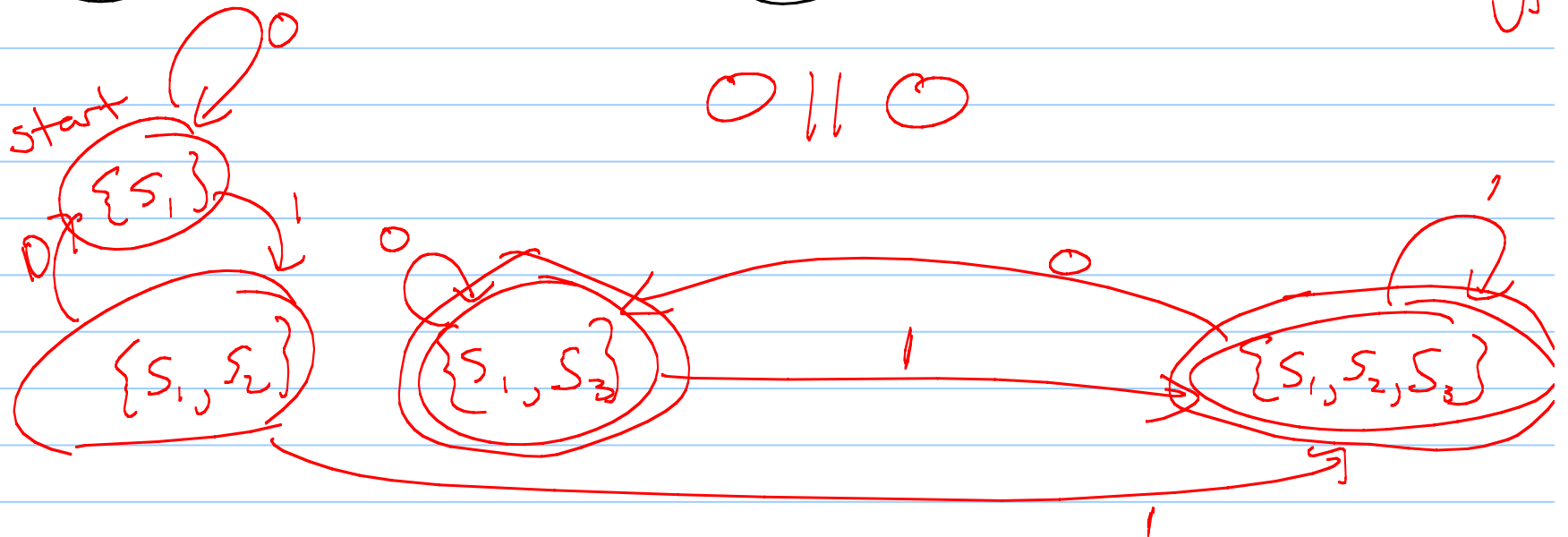NFA has n states
DFA has up to $2^n$ states

# Converting NFAs to DFAs

$\{S_1, S_2, S_3\}$



$\{w \mid w \text{ contains } 11 \text{ as a substring}\}$

$O \, 11 \, O$

# Context free Grammars (+BNF)

Ex:

$$expr \longrightarrow expr \; op \; expr \; |$$
$$(expr) \; | \; -expr \; | \;$$

terminals $\longrightarrow$ id | number

variable

int/float

$$op \longrightarrow + \; | \; - \; | \; * \; | \; /$$

terminals

$$x = 5$$
$$y = 2$$

A **derivation**:   derive  slope * x + intercept

variables (ids)

$expr \Rightarrow expr \ op \ expr$

$\Rightarrow expr \ + \ expr$

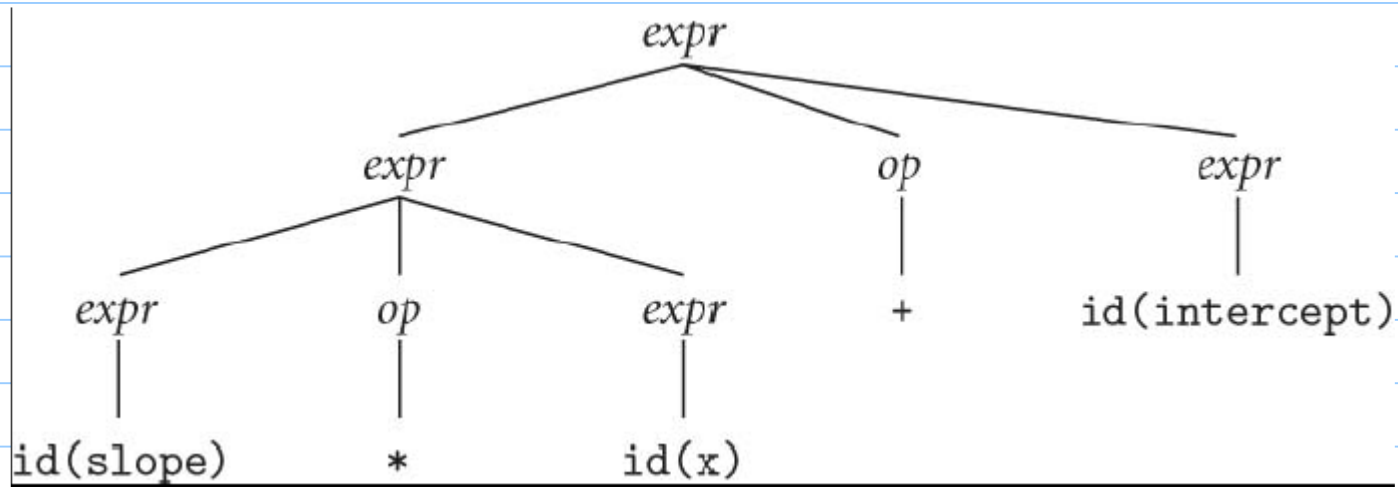$\Rightarrow expr \ + \ id(intercept)$

$\Rightarrow expr \ op \ expr \ + \ id$

$\Rightarrow expr \ * \ expr \ + \ id$

$\Rightarrow expr \ * \ id(x) \ + \ id$

$\Rightarrow id(slope) \ * \ id \ + \ id$

# Derivation tree



*expr*

   *expr*       *op*      *expr*

  *expr*  *op*  *expr*    +   id(intercept)

id(slope)  *  id(x)

(rightmost derivation)

# Ambiguous grammars

expr
├── expr
│   ├── expr — id(slope)
│   ├── op — *
│   └── expr — id(x)
├── op — +
└── expr — id(intercept)

**leftmost derivation** →

expr
├── expr — id(slope)
├── op — *
└── expr
    ├── expr — id(x)
    ├── op — +
    └── expr — id(intercept)

# Grammars

There are infinitely many ways to make a grammar for any context free language.

Problem in the pasing stage: which is better?

(Try to define unambignous grammars.)

# Another example    (from last time)

## Expression grammars : Simple calculator

expr → term | expr add_op term

term → factor | term mult_op factor

factor → id | number | -factor | (expr)

           variable         terminals

add_op → + | -

mult_op → * | /

# Parse Tree

Ex: 3 + 4 * 5

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

# Scanners

Find the syntax (not semantics) of code.

Output tokens.

## 3 types

- Ad-hoc

- Finite automata
  - ~nested case statements
  - ~table & driver

# Ad-hoc   (last time)

If current ∈ { "(", ")", "+", "-", "*" }
     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
etc.

# Ad-hoc approach

Advantage:
    code is fast & compact
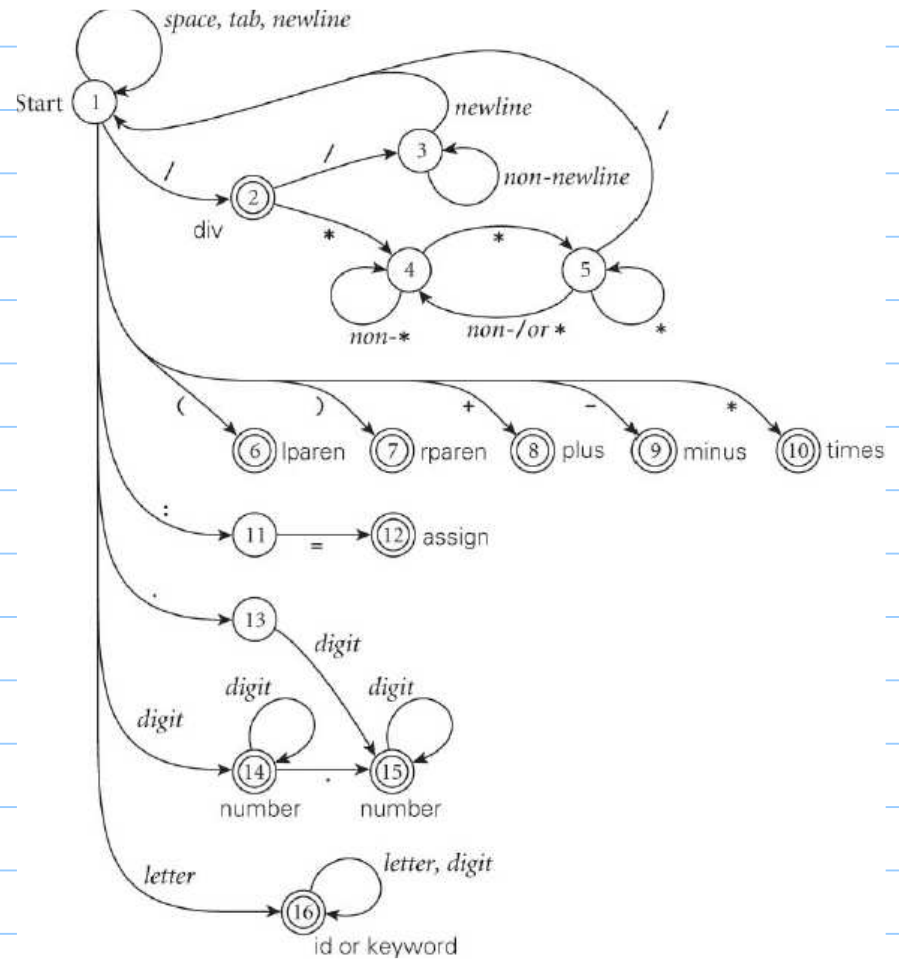
Disadvantage:
    very ad-hoc!
      - hard to debug
       - no explicit representation

# DFA approach

Recall our simple calculator language.
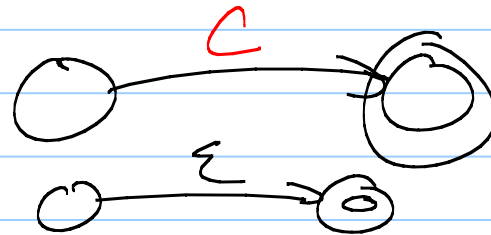
But how to get this DFA, & then how to actually model it?



Start ① space, tab, newline

② div / ③ newline, non-newline

④ * ⑤ non-*, non-/or *, *

( ⑥ lparen ) ⑦ rparen + ⑧ plus - ⑨ minus * ⑩ times

: ⑪ = ⑫ assign

⑬ digit

digit ⑭ number . digit ⑮ number

letter ⑯ letter, digit id or keyword

# Constructing a DFA

Given a regular expression, we can construct an NFA.

Simple NFA:

or

(Base case)

# 3 operations

## Concatenation:

NFA for A
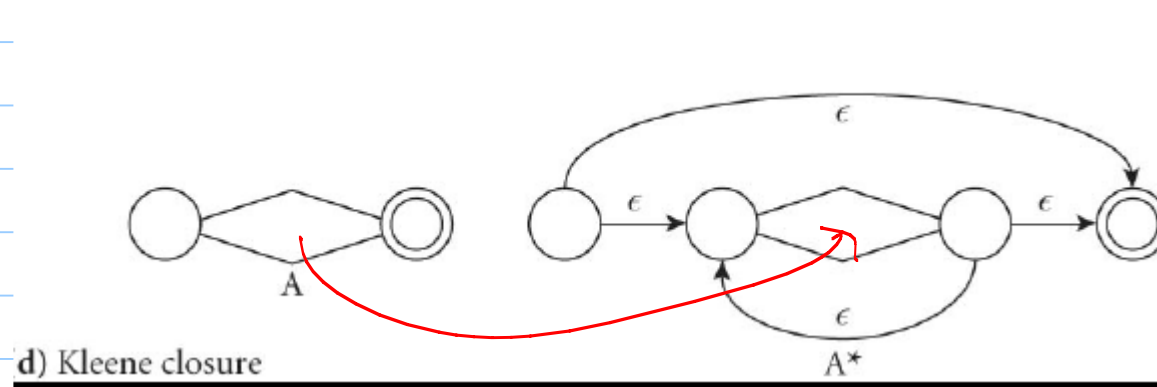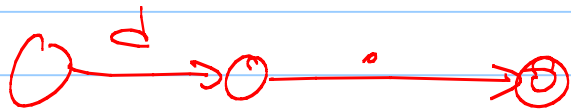
NFA for B

concatenation

AB

both end of A's & start of B's

## Or:

A

B

A|B

both end of A's & start of B's

$\epsilon$  $\epsilon$  $\epsilon$  $\epsilon$

and Kleene closure ($*$).



d) Kleene closure

# Example: decimals $d^*(.d \mid d.) d^*$

## Base:



$[0-9]$

$d^#:$

# Final product:



A | B | C

**Next:** Convert to DFA.
(Lots of states, but same principle as we saw earlier.)

**Result:**

Start
A[1, 2, 4, 5, 8] →d→ B[2, 3, 4, 5, 8, 9] ↻ d
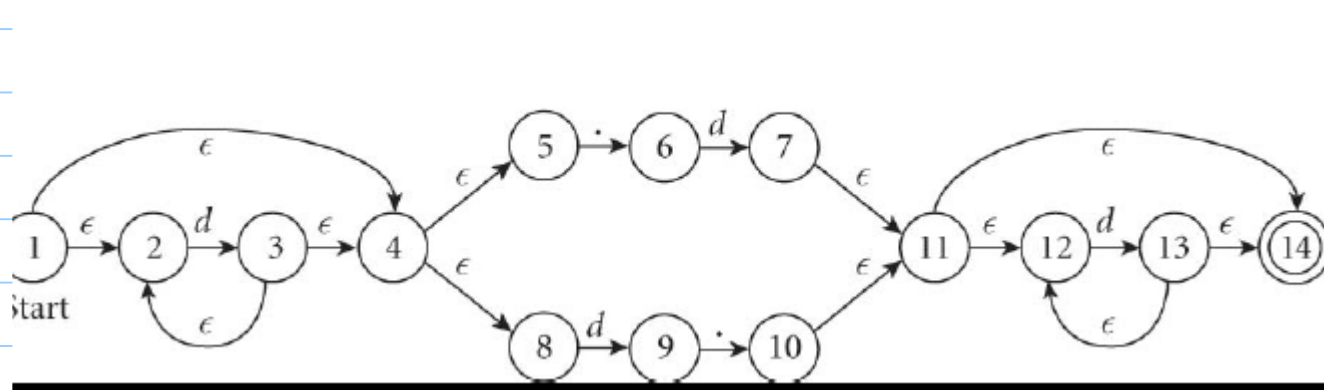
A[1, 2, 4, 5, 8] → C[6]

B[2, 3, 4, 5, 8, 9] → D[6, 10, 11, 12, 14]

C[6] →d→ E[7, 11, 12, 14]

D[6, 10, 11, 12, 14] →d→ F[7, 11, 12, 13, 14]

E[7, 11, 12, 14] →d→ G[11, 12, 13, 14]

F[7, 11, 12, 13, 14] →d→ G[11, 12, 13, 14]

G[11, 12, 13, 14] ↻ d

Note: This
DFA is
a bit redundant.

Not minimal.

Can easily
find the
equivalence
classes.
and minimize.

Start

$A[1, 2, 4, 5, 8]$ $\xrightarrow{d}$ $B[2, 3, 4, 5, 8, 9]$ $\circlearrowright d$

$A$ $\xrightarrow{d}$ $C[6]$

$B$ $\xrightarrow{d}$ $D[6, 10, 11, 12, 14]$

$C$ $\xrightarrow{d}$ $E[7, 11, 12, 14]$

$D$ $\xrightarrow{d}$ $F[7, 11, 12, 13, 14]$

$E$ $\xrightarrow{d}$ $G[11, 12, 13, 14]$ $\circlearrowright d$

$F$ $\xrightarrow{d}$ $G$

# Process to minimize



(a) Start
ABC $\xrightarrow{d,.}$ DEFG

(b) Start
AB → C, AB → DEFG, C $\xrightarrow{d}$ DEFG

(c) Start
A $\xrightarrow{d}$ B, A → C, B → DEFG, C $\xrightarrow{d}$ DEFG
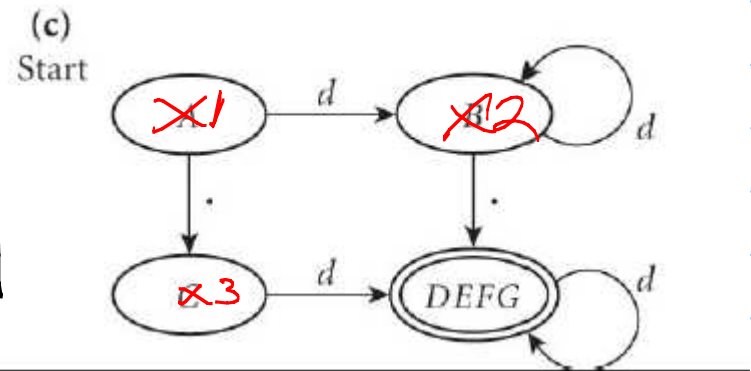
## Now:

Given DFA, generate case statements to simulate it.

```
State = 1
repeat;
    read curr_char
    Case state is:

    1 :   case curr_char = d
              state = 2
          case curr_char = .
              state = 3

    2 :
```



(c)

Start

X1 →(d)→ X2 ↻(d)

X1 →(.)→ X3

X2 →(.)→ DEFG ↻(d)

X3 →(d)→ DEFG

# Scanner Tools

In reality, this DFA is often done() automatically.

Specify the rules of regular language, & the program does this for you.

Many such examples:

Lex (flex), Jlex /Jflex, Quex, Ragel, ...

## Next time:

Lex/Flex : C-style driver

Look for HW on regular expressions, NFA/DFA, & context free languages

Next programming assignment will use flex.