

CS 344 - Scanning

Note Title

1/24/2012

Announcements

- HW2 is due in 1 week

## Last time: Regular expressions

- A character
  - The empty string,  $\epsilon$
  - 2 regular expressions concatenated
  - 2 regular expressions separated by an or (written  $|$ )
  - A regular expression followed by  $*$   
(Kleene star - 0 or more occurrences)
- } base cases

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

$1(0|1)^*0$

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

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

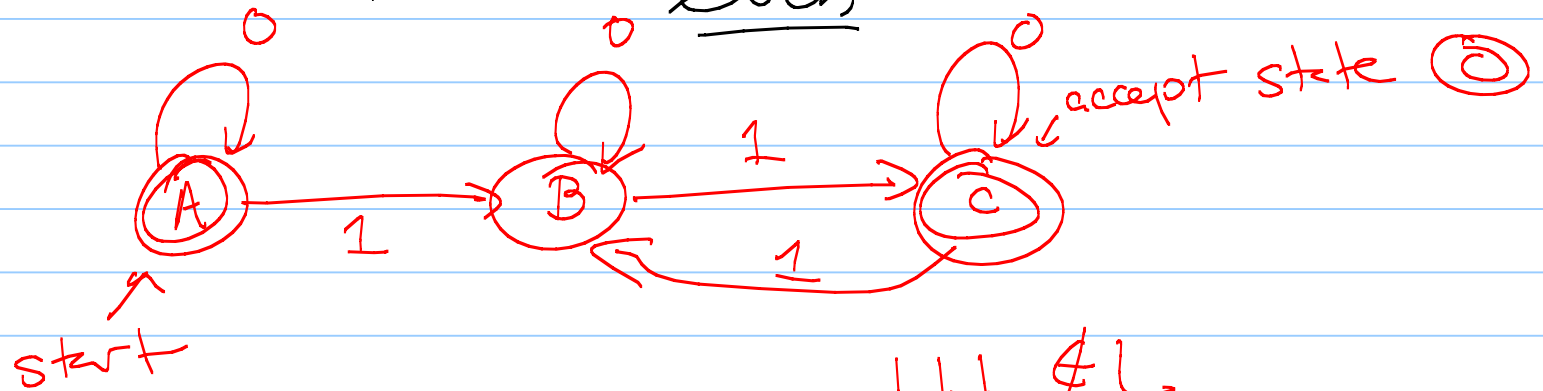
## Deterministic Finite Automate (DFA)

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  $q$  and an input, output a new state

001010 ∈ L

Ex: String of 0's & 1's:  
L = accept if number of 1's is even

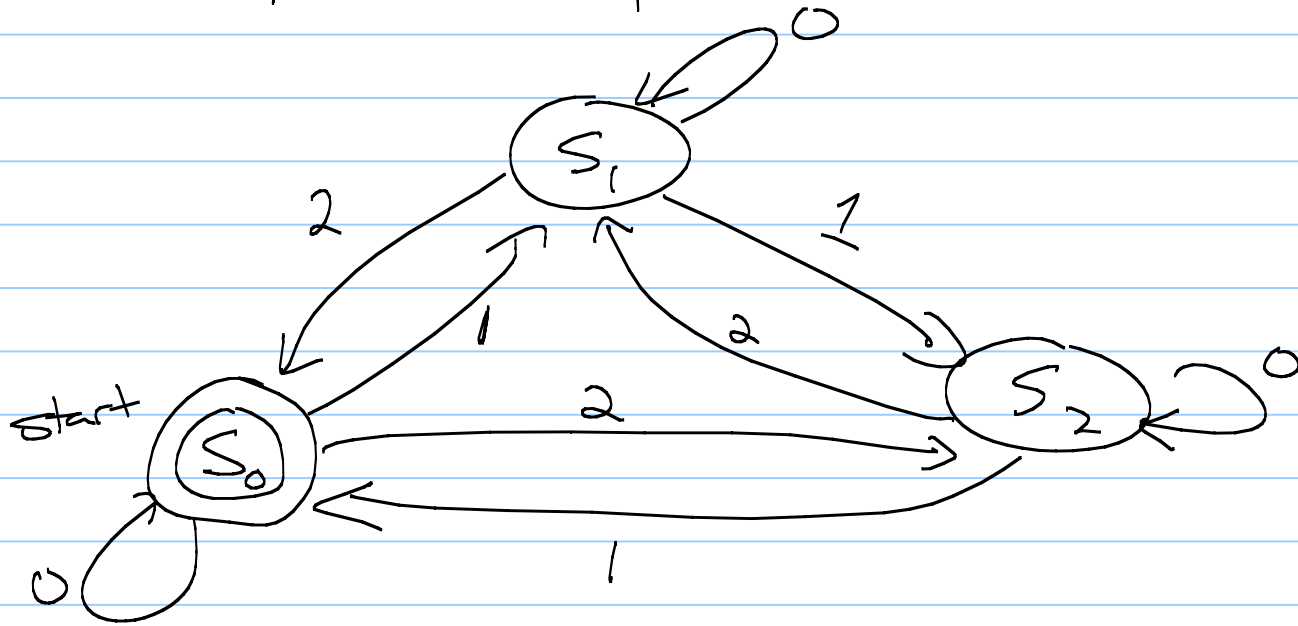


111 ∉ L

1111 ∈ L

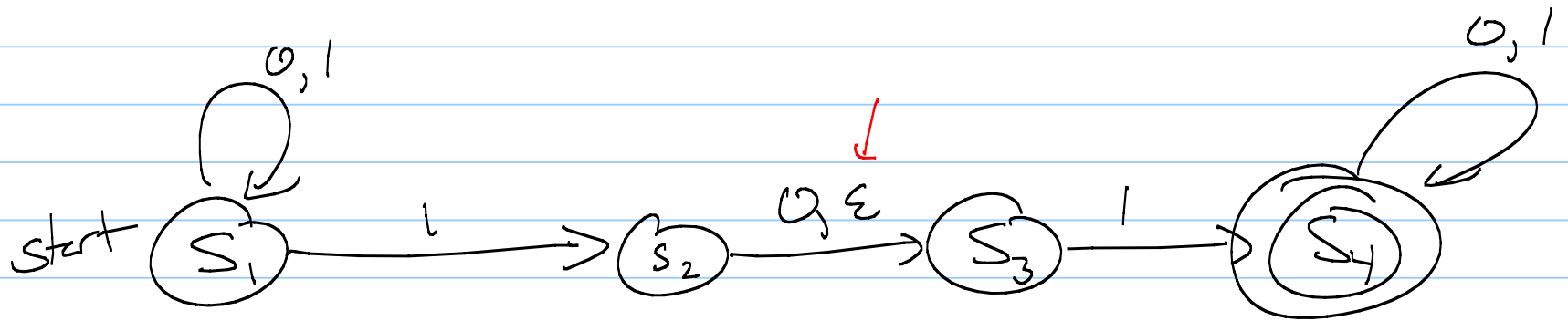
reg exp:  $(0^*10^*10^*)^*$

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



Counts modulo 3  
accepts words  $w \mid \text{sum} \equiv 0 \pmod 3$

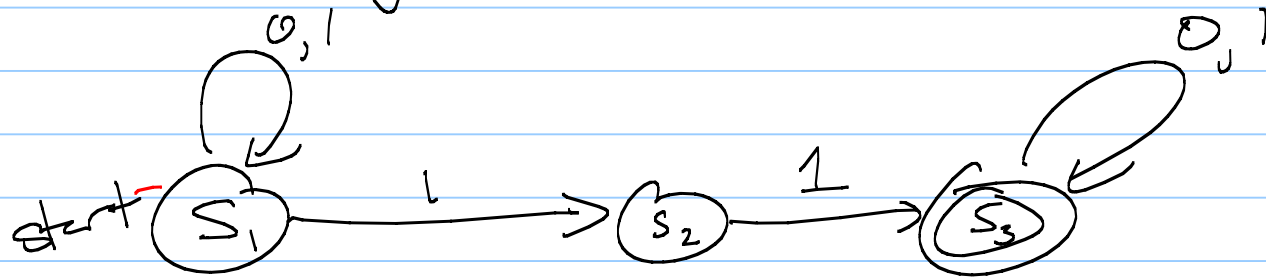
NFAs: DFAs w/ ambiguity



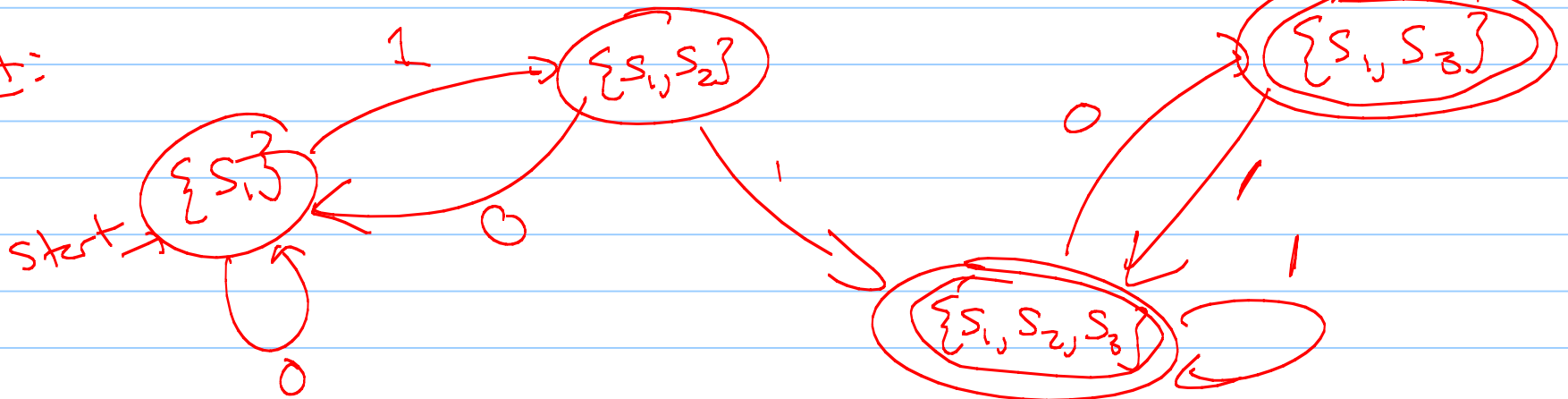
→ 10

Thm: NFAs & DFAs are equivalent.

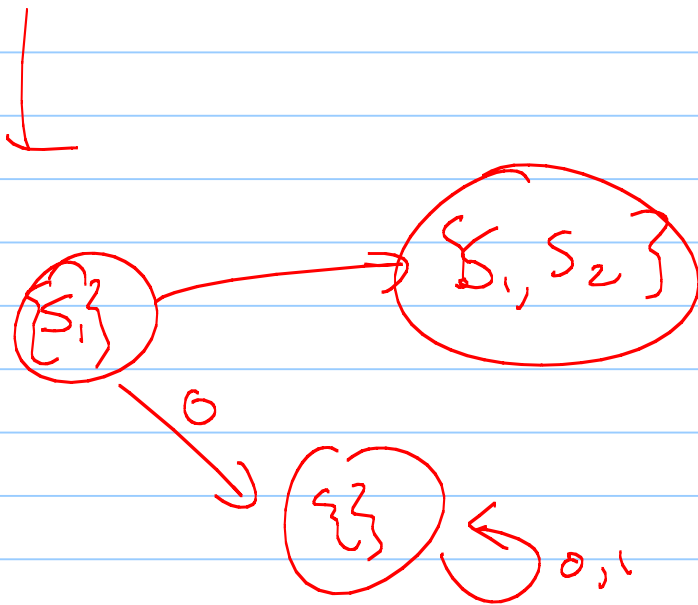
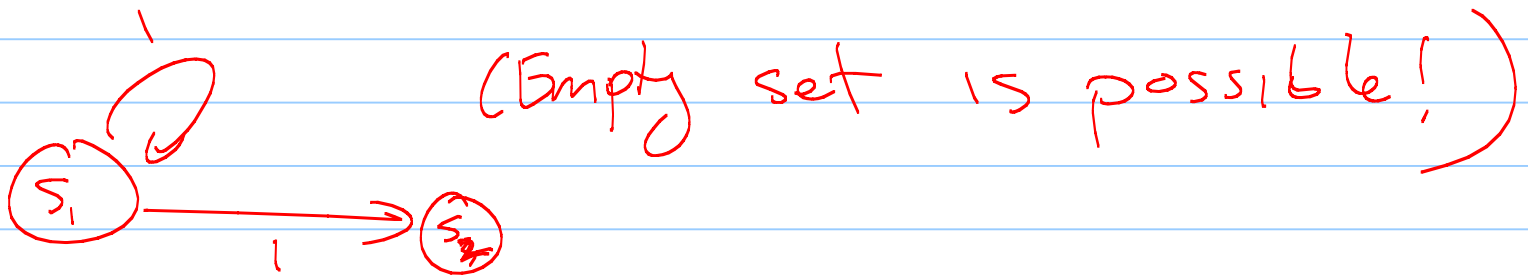
# Converting NFAs to DFAs (p.57)



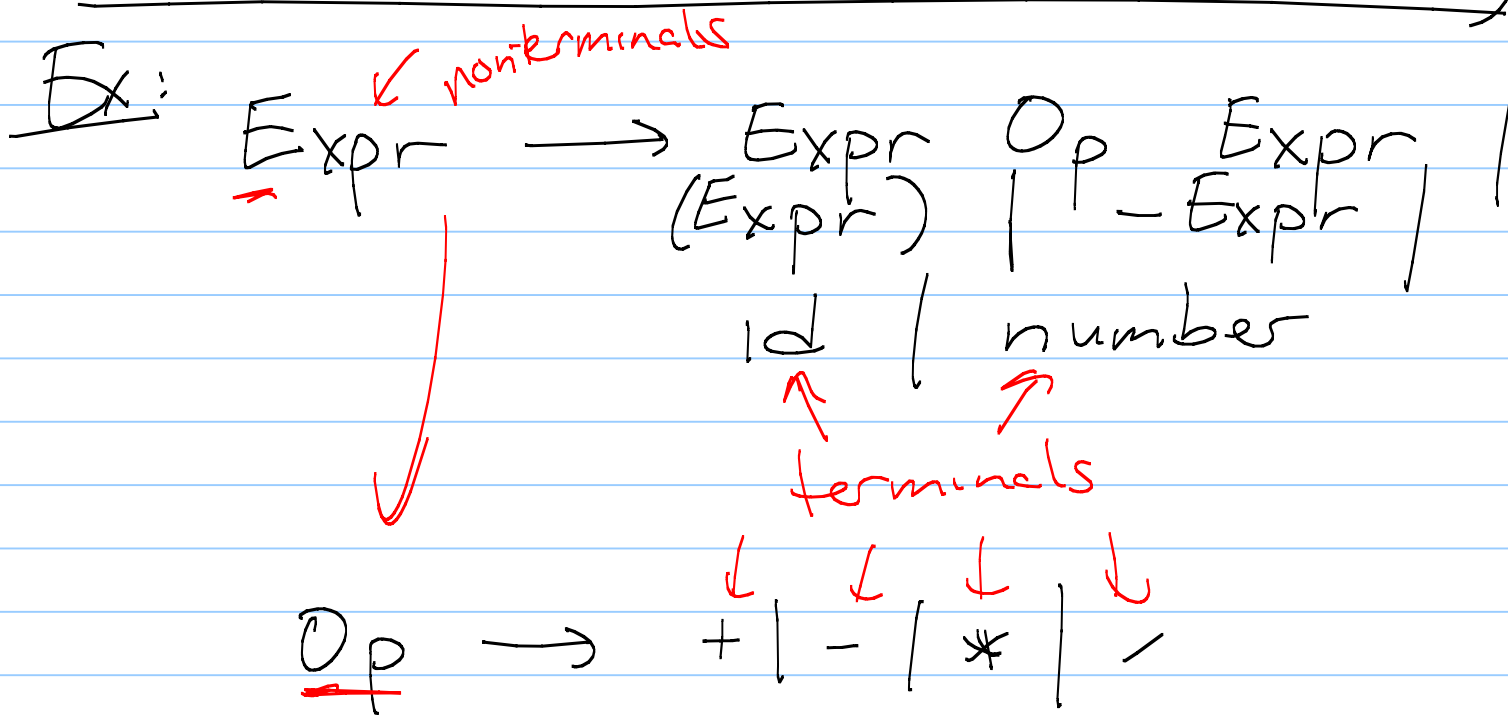
DFA:







# Context free Grammars (BNF)



A derivation: derive  $\frac{1}{2} * x + \text{intercept}$

variables (ids)  $\swarrow$   $\searrow$

$\text{Expr} \Rightarrow \text{Expr} \underline{\text{Op}} \text{Expr}$

$\Rightarrow \text{Expr} + \text{Expr}$

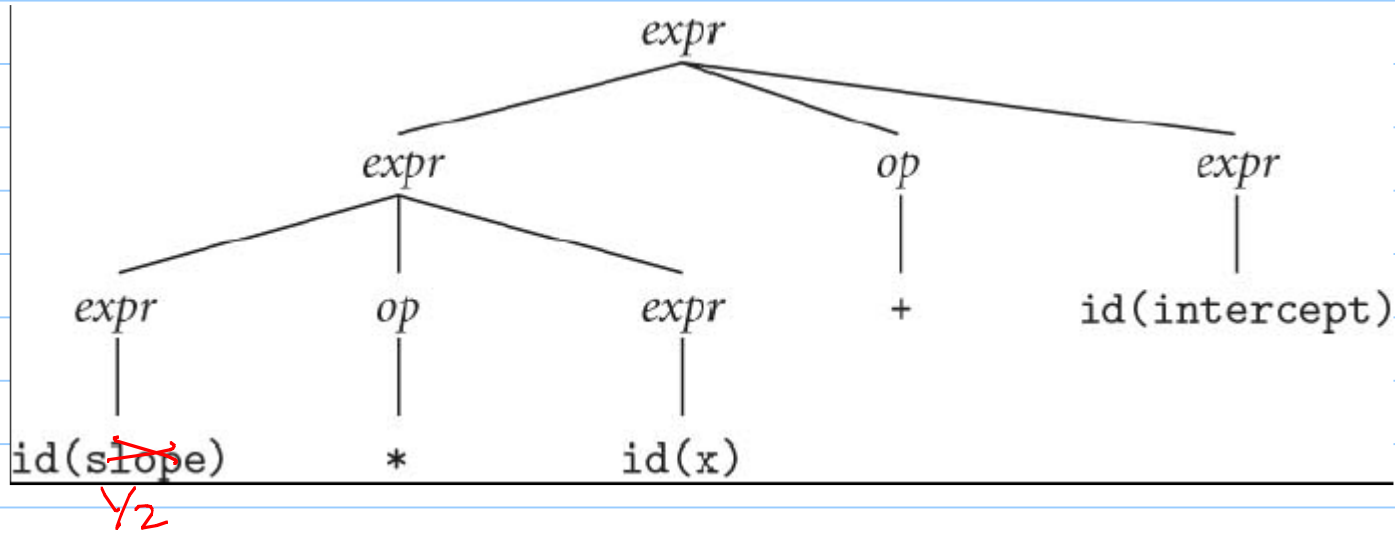
$\Rightarrow \underline{\text{Expr}} + \text{id}(\text{intercept})$

$\Rightarrow \text{Expr} \underline{\text{Op}} \text{Expr} + \text{id}(\text{intercept})$

$\Rightarrow \underline{\text{Expr}} * \underline{\text{Expr}} + \text{id}(\text{intercept})$

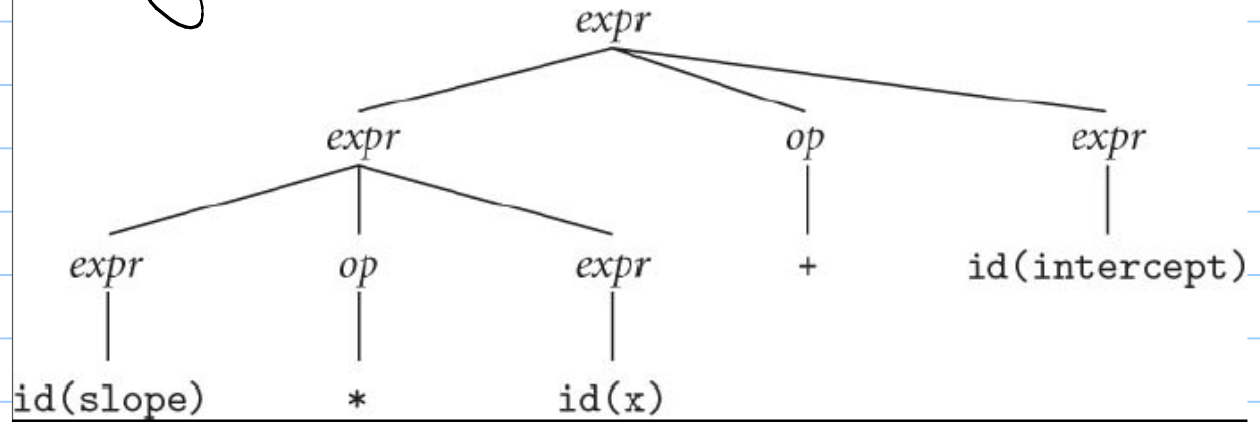
$\Rightarrow \text{number}(\frac{1}{2}) * \text{id}(x) + \text{id}(\text{intercept})$

# Derivation tree

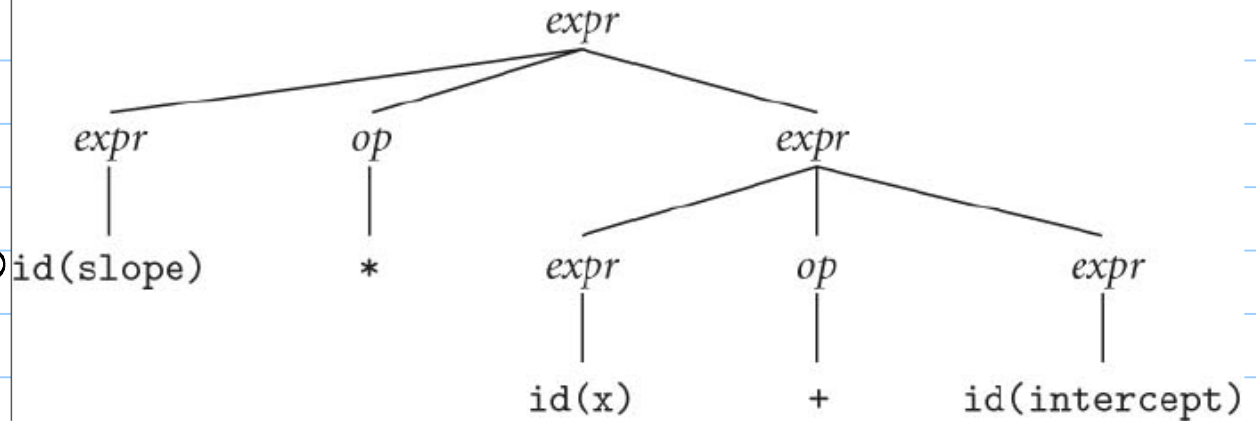


(rightmost derivation)

# Ambiguous grammars



leftmost  
derivation →



# Grammars

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

Problem in the parsing stage:  
which is better?

(Try to define unambiguous grammars.)

Another example (from last time)

goal: avoid ambiguity from last ex.

Expression grammars: simple calculator

Expr  $\rightarrow$  Term / Expr Add\_op Term

Term  $\rightarrow$  Factor | Term Mult\_op Factor

Factor  $\rightarrow$  id | number | -Factor | (Expr)

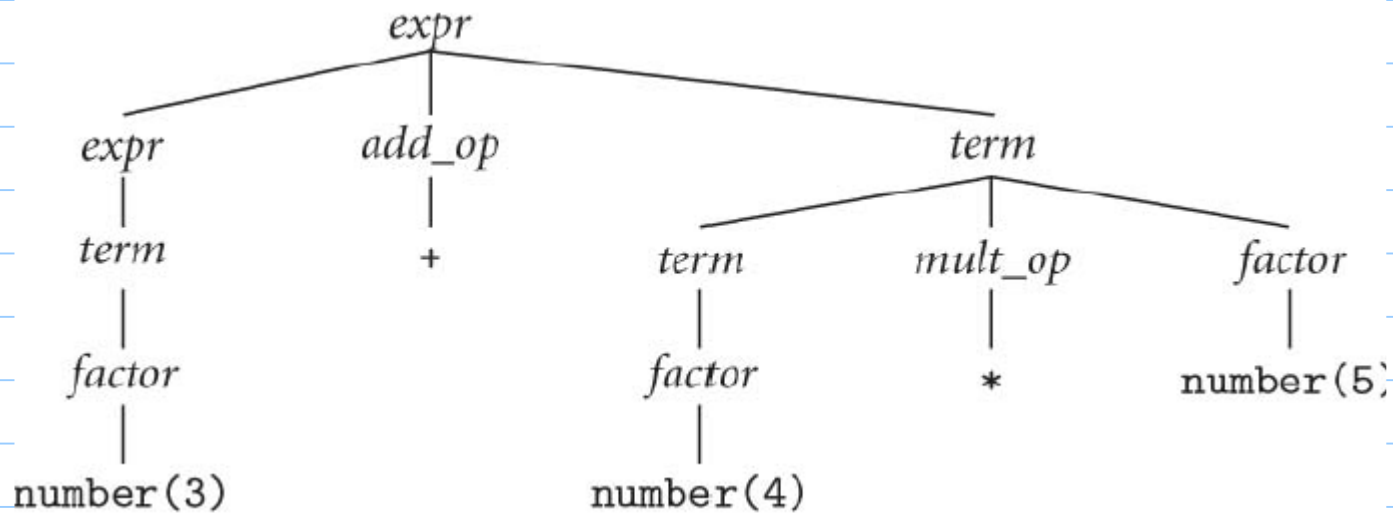
Add\_op  $\rightarrow$  + | -

Mult\_op  $\rightarrow$  \* | /

# Parse Tree

Ex:  $3 + 4 * 5$

(avoids rightmost derivation)





Scanners: do this in code

Find the syntax (not semantics)  
of code.

Output tokens.

A few types:

- Ad-hoc

- Finite automata  $\leftarrow$  DFA/NFA
  - nested case statements
  - table + driver  $\leftarrow$  Simulates DFA

① Ad-hoc : case based code <sup>variable assignment</sup>

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 "/" or "newline" (resp.)

else return divide

etc.

## Ad-hoc approach

Advantage:

code is fast & compact

Disadvantage:

very ad-hoc!

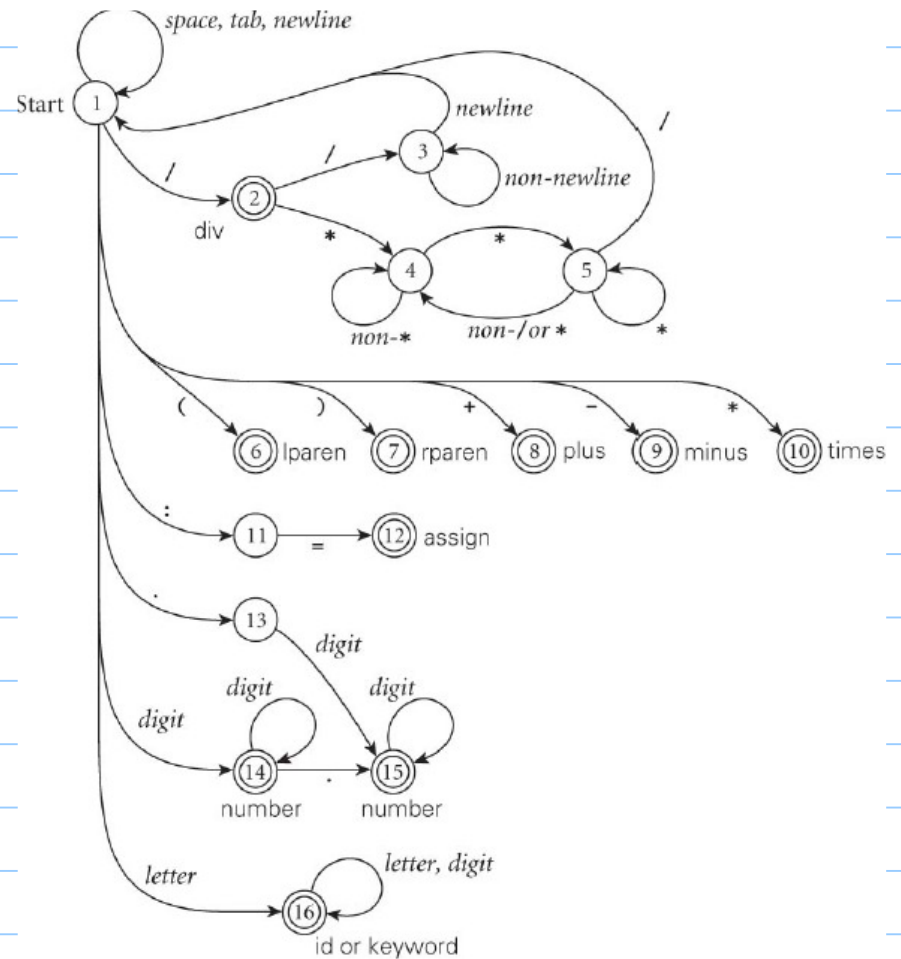
- hard to debug

- no explicit representation

# DFA approach

Recall our simple calculator language.

But how to get this DFA  
or then how to actually model it?

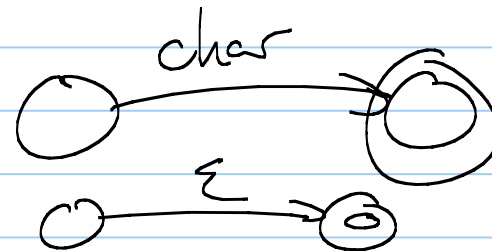


## Constructing a DFA (p57)

Given a regular expression, we can construct an NFA.

Simple NFA:

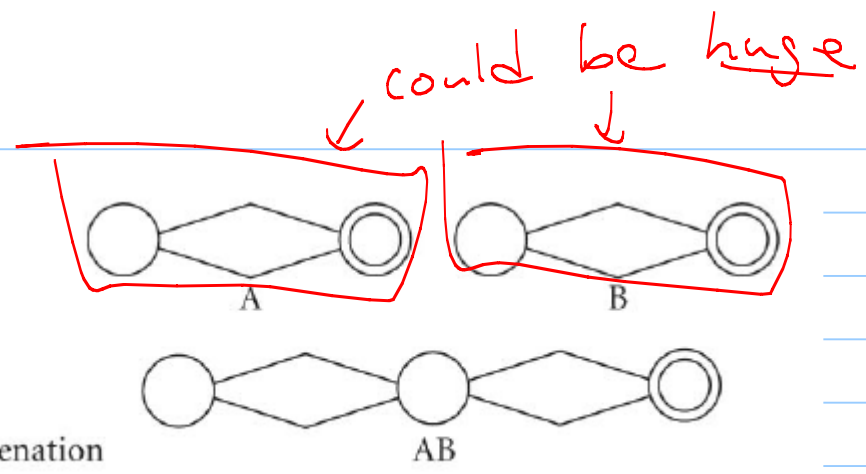
or



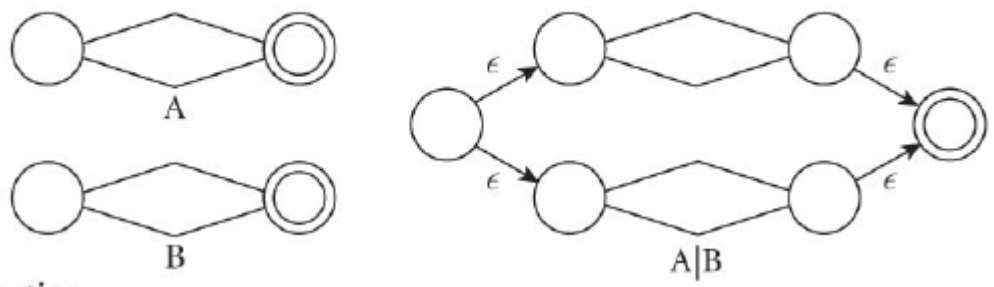
(Base case)

# 3 operations

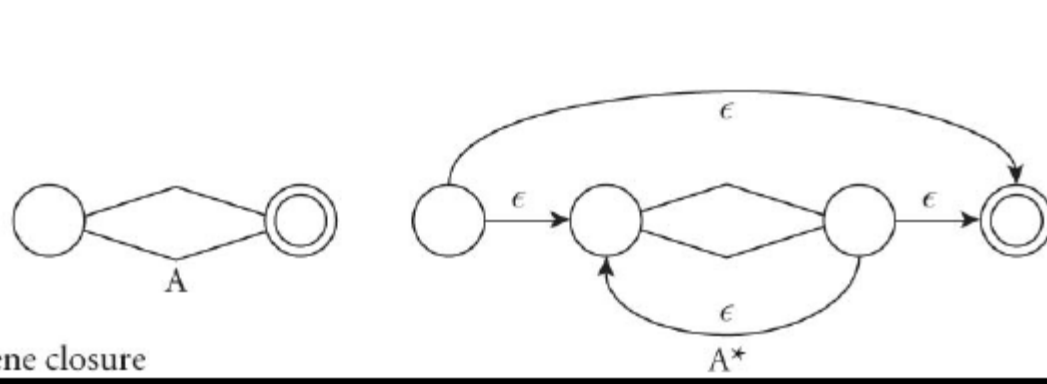
## Concatenation:



Or:

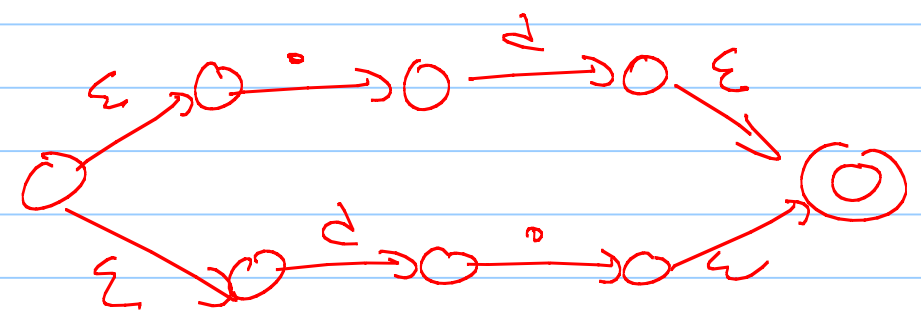
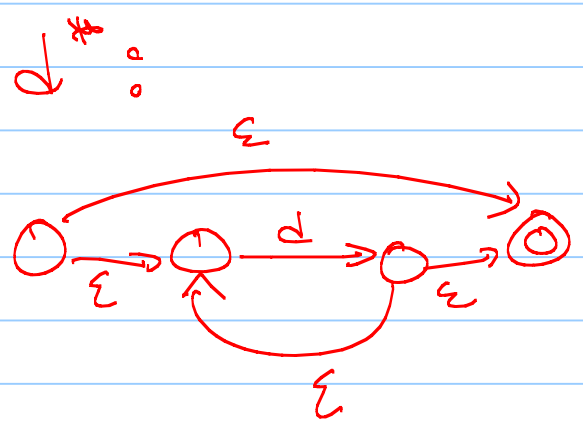


and Kleene closure ( $*$ ).



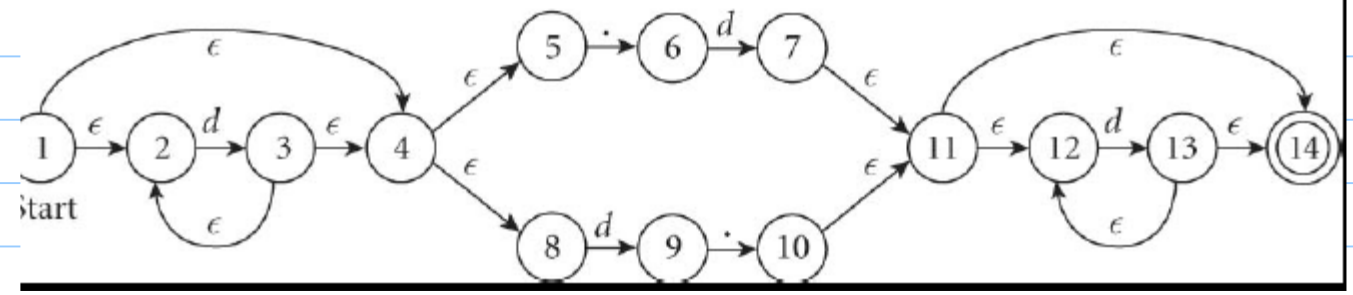
d) Kleene closure

Example: decimals  $d^*(.d|d.)d^*$



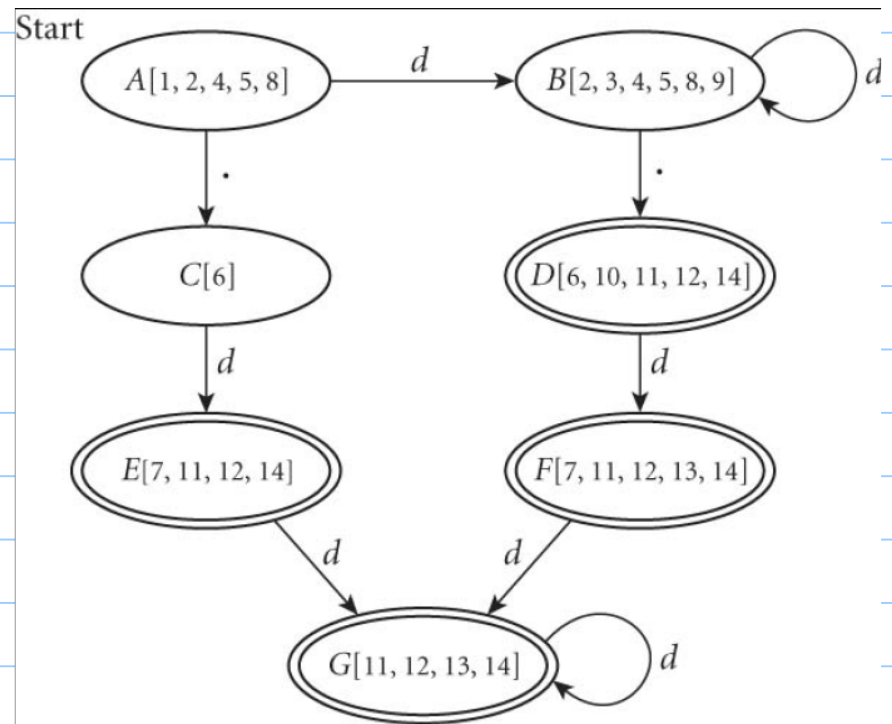


Final product :



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

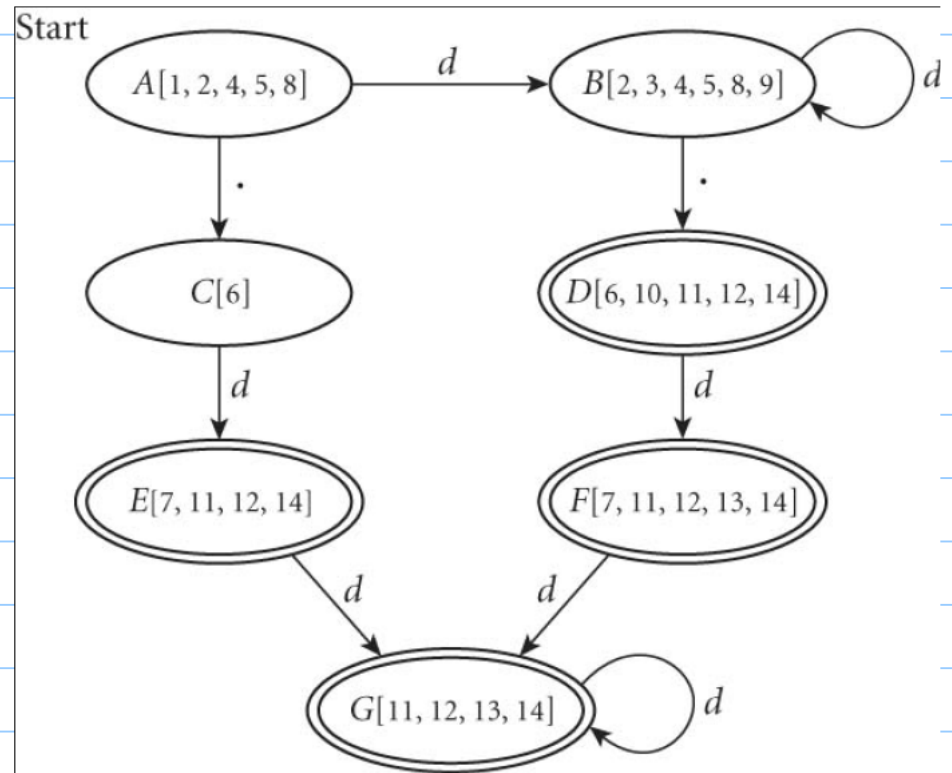
Result ∴  
(see p. 57-58)



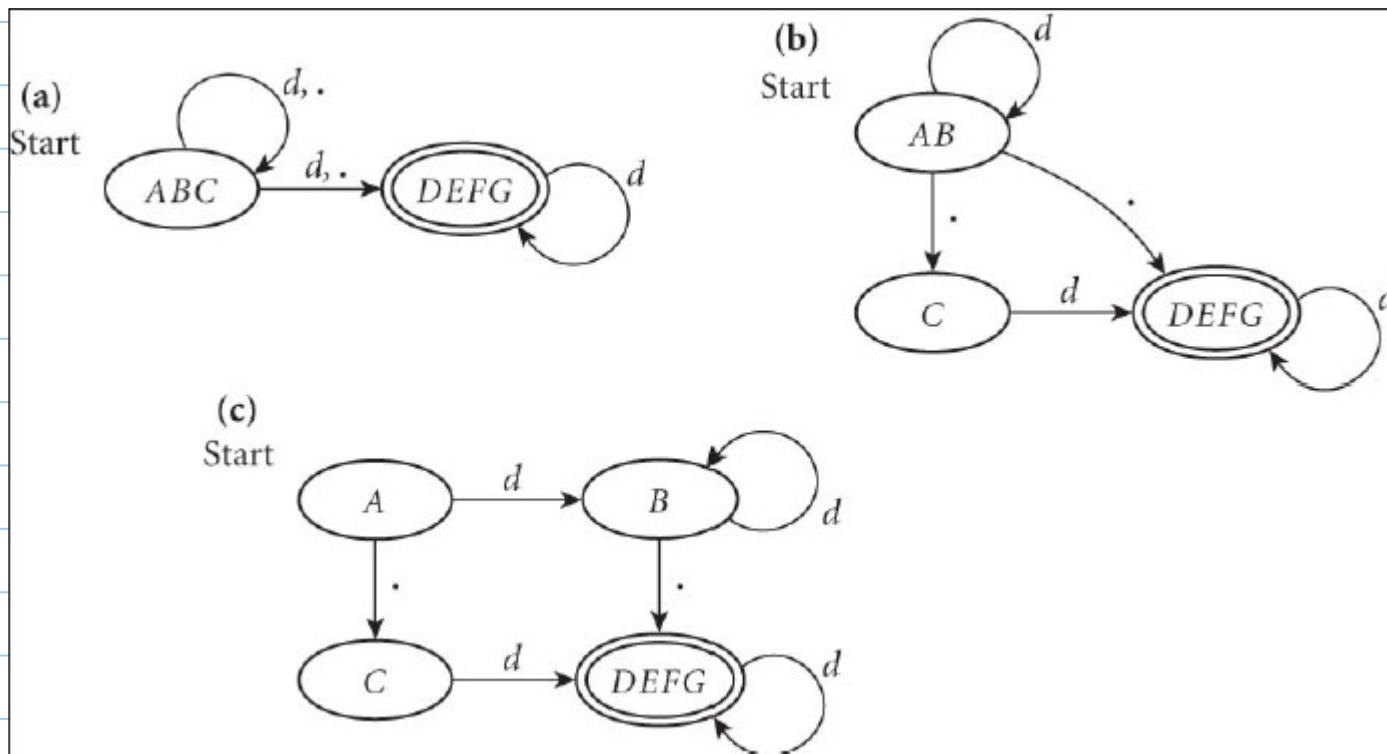
Note: This DFA is a bit redundant.

Not minimal.

Can easily find the equivalence classes and minimize.



# Process to minimize



Now:

Given DFA, generate case statements to simulate it.

State = 1

repeat;

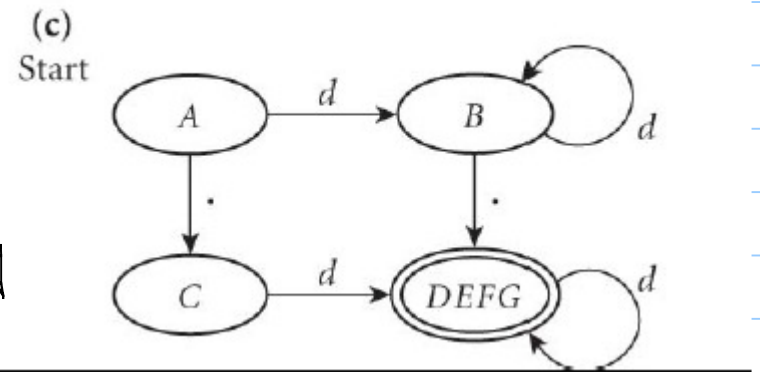
read curr\_char

case state is:

A: case curr\_char = d  
state = B

case curr\_char = .  
state = C

B:



# Scanner Tools

In reality, this DFA is often done automatically.

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

Many such examples:

Lex (flex), Jflex / 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.