

# Lecture 11+12

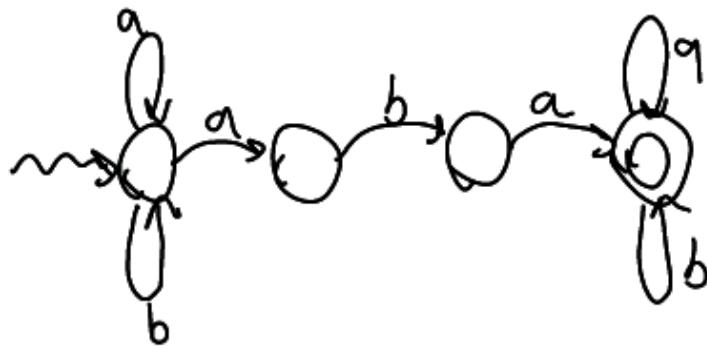
## Regular Expressions

*CS 241: Foundations of Sequential Programs*  
Fall 2014

Troy Vasiga et al  
University of Waterloo

## Review

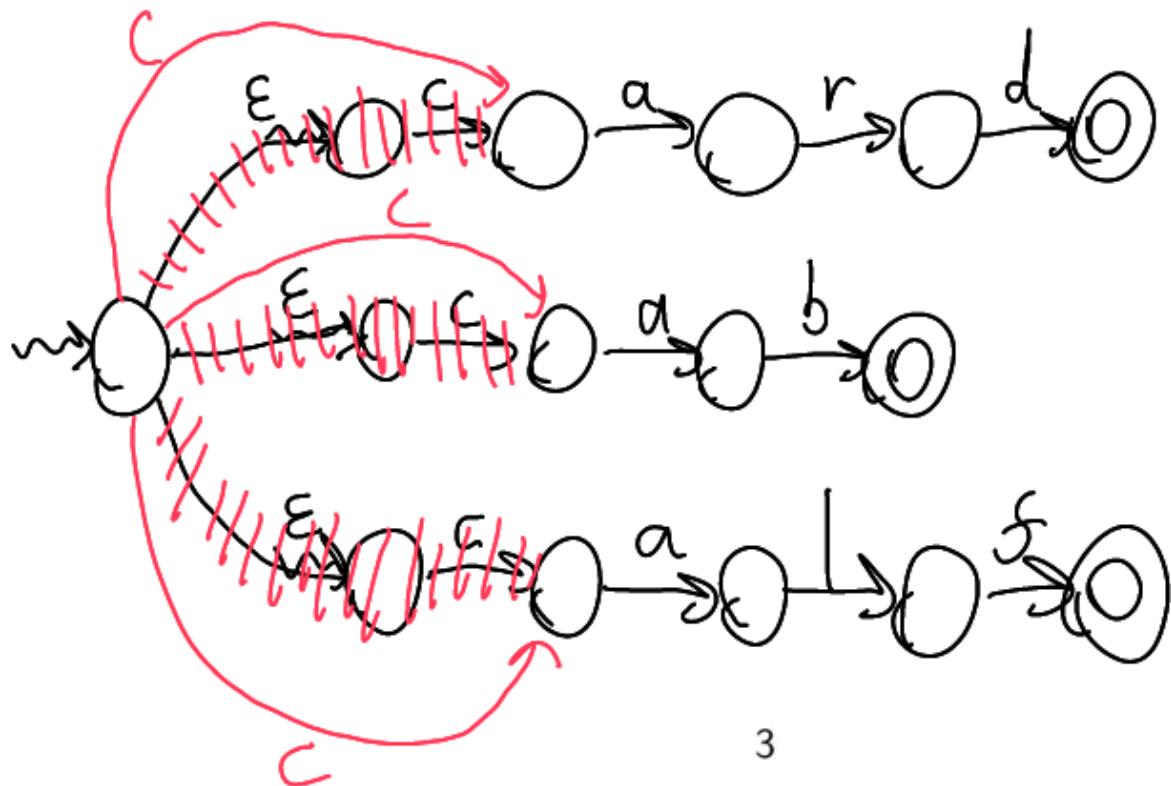
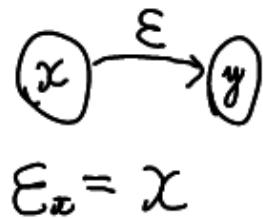
- ▶ An example comparing DFAs and NFAs



- ▶ Can convert between an NFA and DFA using the subset construction
  - ▶ define each set of states that can be occupied at the same step
  - ▶ one state in the DFA for each unique set of states in the NFA

## $\epsilon$ -NFAs

- ▶ allows transition between states on “no input”
- ▶ can be used as “glue” for joining machines together
- ▶ example:  $L = \{ \text{card, cab, calf} \}$
- ▶ it is not surprising that  $\epsilon$ -NFAs can be converted to NFAs



## Regular Expressions

specify a regular language

Defined recursively: a regular expression (RE) is

- ▶  $\emptyset$ , or  $L = \{\}$
  - ▶  $\epsilon$ , or  $L = \{\epsilon\}$
  - ▶  $a$ , where  $a \in \Sigma$
  - ▶  $E_1 E_2$  where  $E_1$  and  $E_2$  are REs
  - ▶  $E_1 | E_2$  where  $E_1$  and  $E_2$  are REs
  - ▶  $E^*$  where  $E$  is a REs
- finite languages ← base cases
- concatenations
- union/alternating "or"
- repetition

zero or more copies

recursive cases

## RE examples

- $L = \{cab, car, card\}$

$cab \mid car \mid card$

or:  $ca(b|r|rd)$   
 $ca(b|r(d|\epsilon))$

- $\Sigma = \{a\}, L = \{w: w \text{ contains even } \#s \text{ of } a's\}$

$(aa)^*$        $\epsilon, aa, aaaa, \dots$

- $\Sigma = \{a, b\}, L = \{w: w \text{ contains even } \#s \text{ of } a's\}$

$b^*(ab^*ab^*)^*$

## More RE examples

- $\Sigma = \{a, b\}$ ,  $L = \{w: \text{contains either } aa \text{ or } bb\}$

$$\begin{array}{c} aa \mid bb \\ (a|b)^* (aa|bb) (a|b)^* \end{array}$$

- $\Sigma = \{a, b\}$ ,  $L = \{w: \text{contains no occurrence of } aa \text{ or } bb\}$

Ex

## (Unix) Shorthands and Extensions

- ▶  $[a - z] \leftarrow (a|b|c|\dots|z)$    
regular expression

- ▶  $E^+ \leftarrow EE^*$   
One or more copies
- ▶ others may occur in your future readings

## RE to $\epsilon$ -NFA

Convert piece by piece. Recall:  
a regular expression (RE) is

- ▶  $\emptyset$ , or



- ▶  $\epsilon$ , or



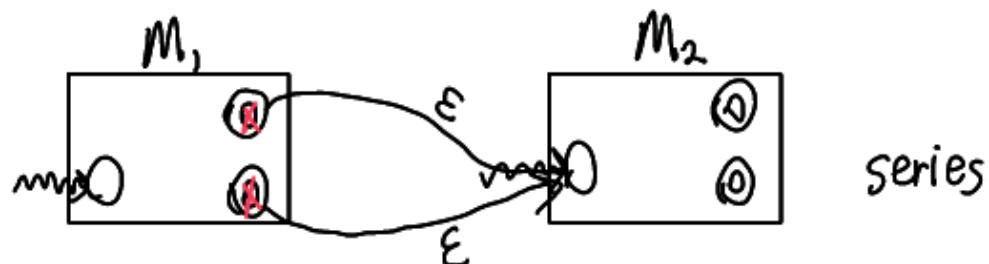
- ▶  $a$ , where  $a \in \Sigma$



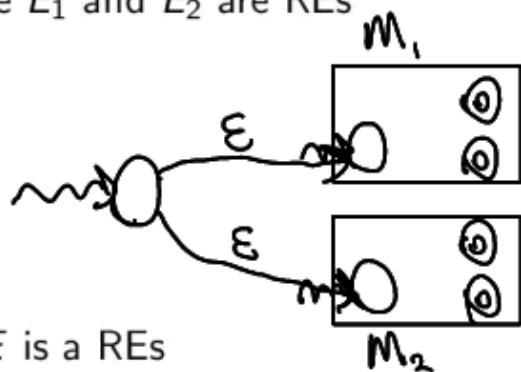
(continued on next slide)

## RE to $\epsilon$ -NFA (continued)

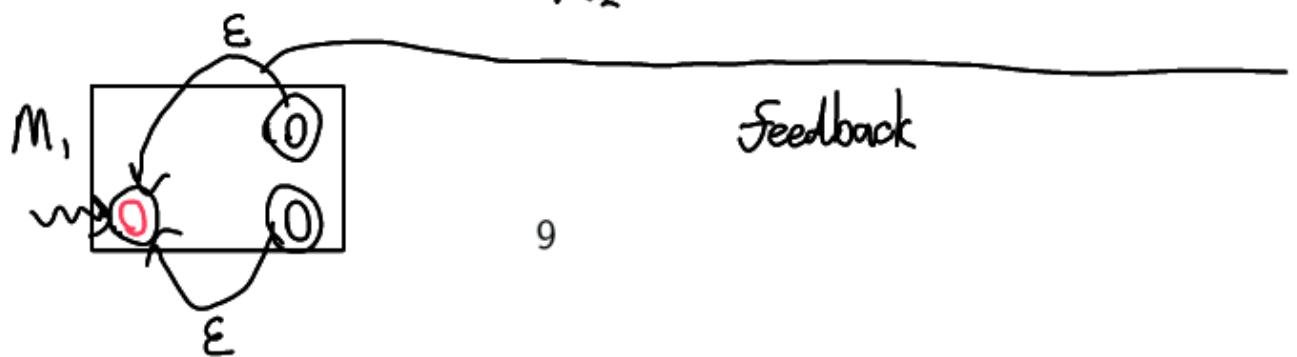
- ▶  $E_1 E_2$  where  $E_1$  and  $E_2$  are REs



- ▶  $E_1 | E_2$  where  $E_1$  and  $E_2$  are REs



- ▶  $E^*$  where  $E$  is a REs

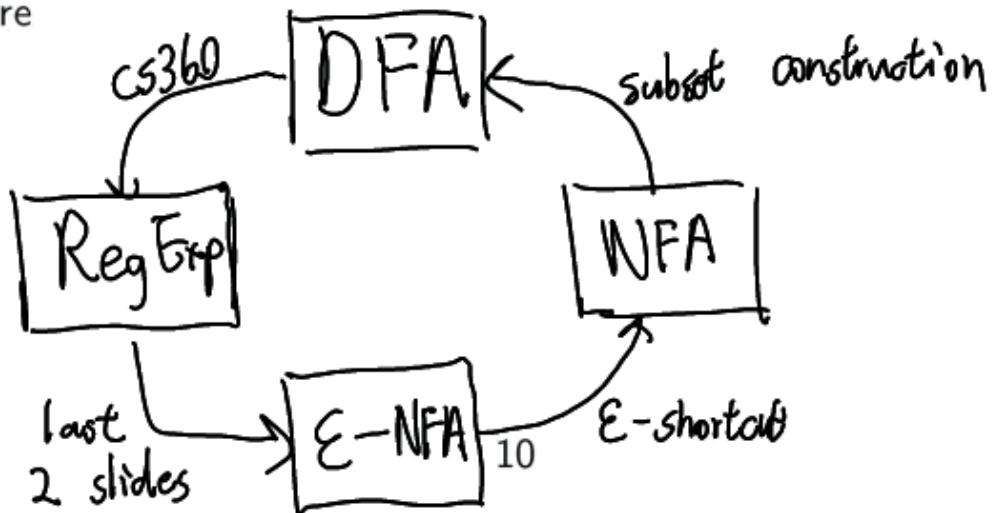


## Circle of Life

Definition of Regular Language:

specified by a regular expression  
recognized by a DFA  
NFA  
 $\epsilon$ -NFA

A picture



## Practical Applications of DFAs

- ▶ Most real-world examples do not care about recognizers (DNA match may be the exception)
- ▶ Mostly, DFAs are used for:
  - ▶ transforming/transducing input
  - ▶ searching in text
  - ▶ scanning/translating

yes/no  
↑

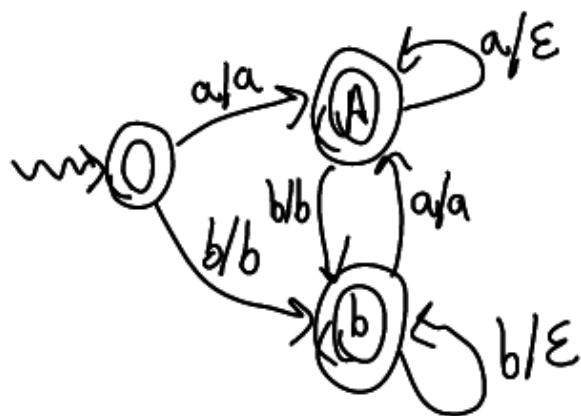
## Transducers

A transducer is a DFA with output. That is, transitions look like input/output.

Example 1: Remove stutters from  $\Sigma = \{a, b\}$ .

$aabbbaaa \Rightarrow aba$

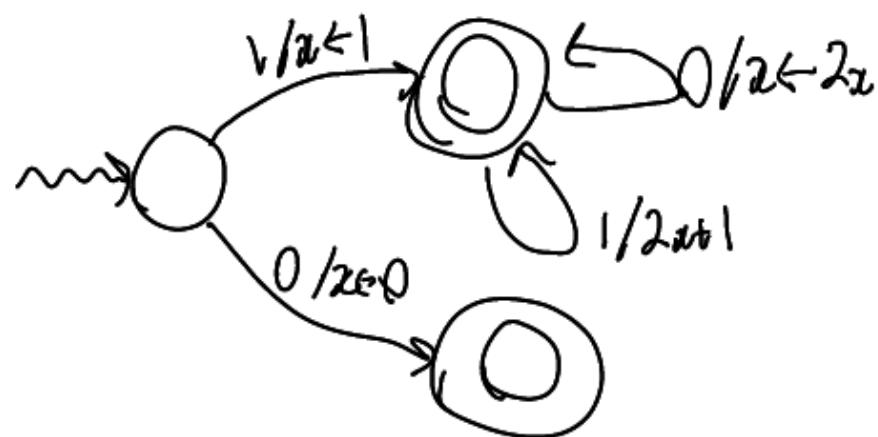
$baabbb \Rightarrow bab$



## Transducers

Example 2: Binary integers from  $\Sigma = \{0, 1\}$ .

$0 \rightarrow 0$       Calculate decimal value  
 $1010 \rightarrow 10$       into the variable  $x$ .  
 $10000 \rightarrow 16$



## Final word about transducers

- ▶ Mealy Machine
- ▶ Moore Machine
- ▶ “translating” does not “give meaning”

## Searching

Goal: Find the first occurrence of a string  $p \in L$  inside a (larger) string  $T$ .  
Solution (naive):

- ▶  $\underline{aab}\leftarrow P$
  - ▶  $T\leftarrow aacubbaabahaca\bar{a}a$
  - ▶ Problems:
- 

Knuth  
Morrist  
Pratt

CS240

This turns out to be difficult, but does rely on DFAs (see CS240).

DFA for MIPS

lexer

Contextual DFA

- ▶ It is easy to write a DFA to recognize an individual token type
- ▶ It is not difficult to combine these into one DFA that recognizes a word as some token type
- ▶ See the DFA for MIPS on the CS241 website

## Scanning

The scanning problem:

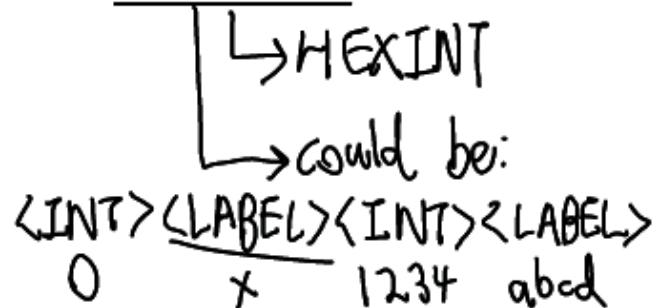
**Input:** some string  $w$  and a language  $L$

"add \$1, \$2, \$3"

**Output:**  $w_1 w_2 \dots w_n = w$  where  $w_i \in L$  for all  $i$

<ADD> <REG> <COMMA> <REG> ...

There may be more than one possible answer: 0x1234abcd



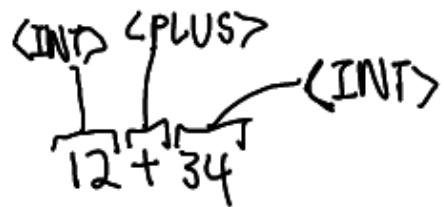
- We solve this problem by...

repeatedly runs the dfa

## Simplified Maximal Munch

Input:  $c_0 c_1 c_2 \dots c_{k-1}$

```
i = 0
state = START
loop
    newstate = ERROR ← assume worst case
    if ( i < k ):
        newstate = δ(state,  $c_i$ ) ← If not at the end of the
        if newstate == ERROR: string, look up next state
            if state is not a final state:
                report an error and exit
            if state is not WHITESPACE:
                output appropriate token
            state = START ← restart again!
            if i == k:
                exit
        else:
            state = newstate
            i = i + 1
    ] move ahead in my input
```



## Next Steps

