

The "traditional" compiler workflow presented in class creates an intermediate format in the middle-end. Why not translate directly from the AST (abstract syntax tree) to final code?



### **Project 1**

- Out tonight
- Add Flex rules for our language
- Might want to find a project partner!



What should we call our language?

## Surveys (Mostly) Processed Housekeeping

- We will do flipped Wednesdays
- I'll try to put out some 510 review

KU | EECS 665 | Drew Davidson

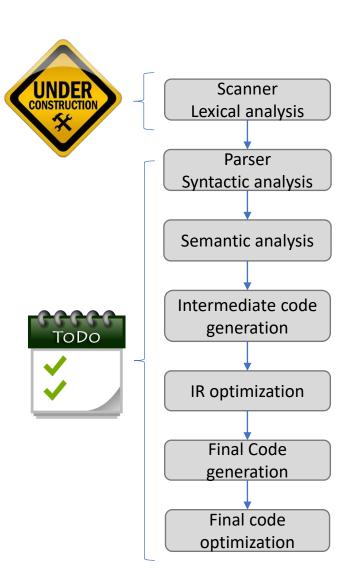
# CONSTRUCTION

2 – Implementing Scanners

## Compiler Construction Progress Pics

## **Currently working on lexical analysis concepts**

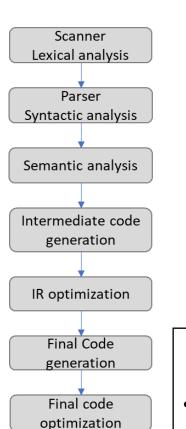
 Convert the character stream from the user into the token stream (the "words" of the programming language)



### Last Time

Review: Overview

Introduced a definition of a compiler and its workflow



Used RegExs to *specify* languages of tokens

<u>Token</u>

Integer Literal

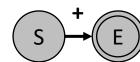
**RegEx** 

0|(1|2|...|9)(0|1|...|9)\*

Claimed RegExs can be automatically translated to *recognize* languages

RegEx

**DFA** 



In lab we showed a tool that automatically does this translation (Flex)

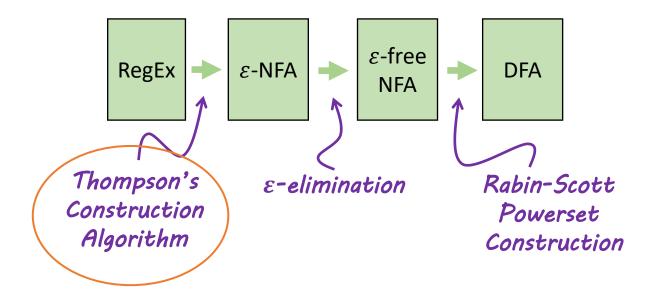
### **You Should Know**

- The phases of the compiler and what each step does
- What errors the various phases catch
- How to specify a token's lexemes with a regex

### Lecture Overview

Lecture 2 – Implementing Scanners

### Walk through the translation process formally

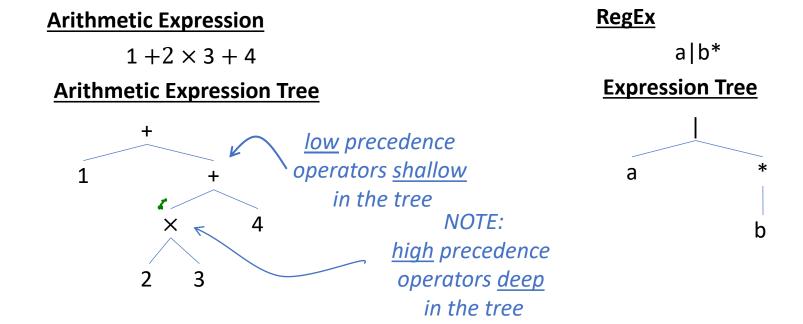


## Key Concept Thompson's Construction Algorithm

### Use an expression tree:

- Leaf: atomic operand
- Branch: operations joining subtrees

### **Expression Tree Examples**



## Thompson's Construction Intuition Thompson's Construction Algorithm

### **Two-Step Process:**

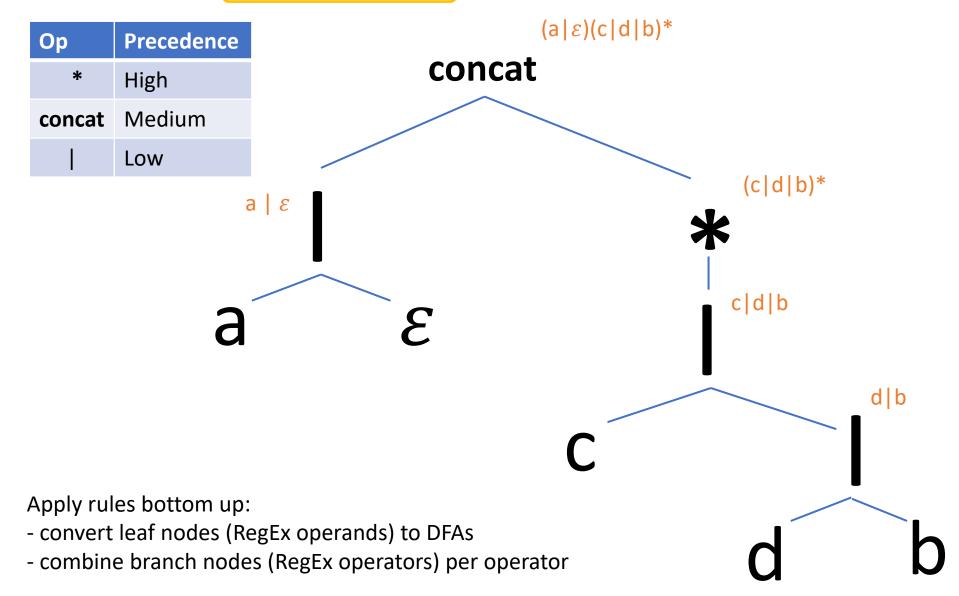
- Break the RegEx down to the simplest units with "obvious" FSMs (i.e. expression tree leaves)
- Combine the sub-FSMs according to operator rules (i.e. expression tree branch rules)

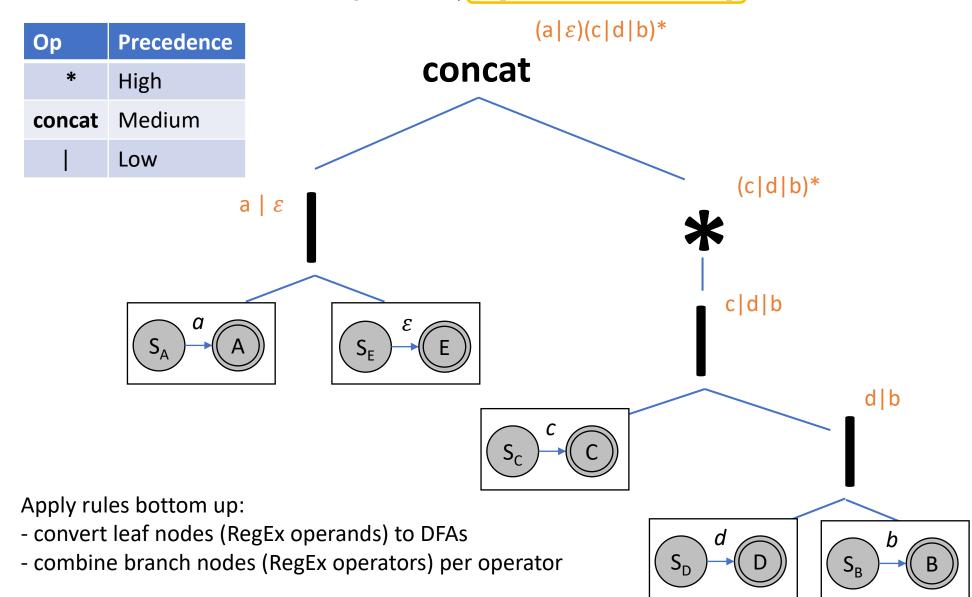


Recombobulate the Regex into an FSM piecewise

### Thompson's Construction Alg.

Build the RegEx Tree | Replace nodes bottom-up



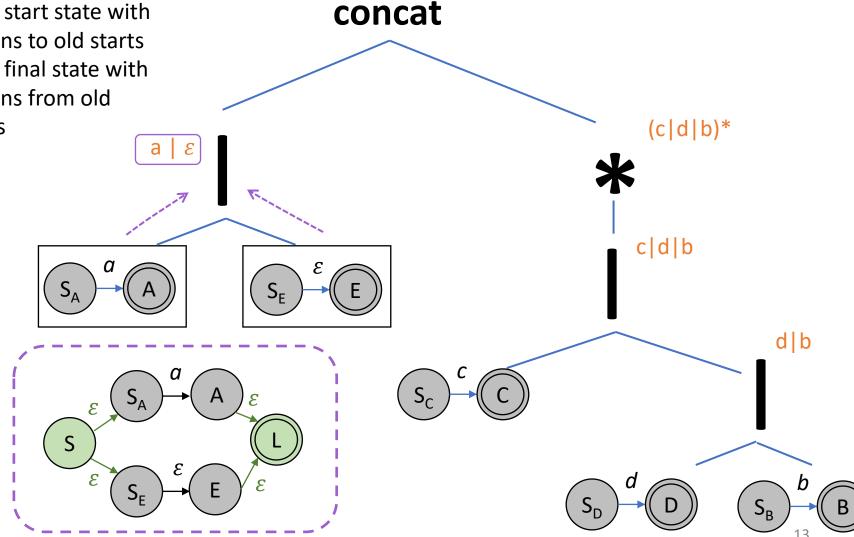


 $(a|\varepsilon)(c|d|b)^*$ 

### Alternation:

New start state with  $\varepsilon$ -trans to old starts

New final state with  $\varepsilon$ -trans from old finals



### Alternation:

New start state with  $\varepsilon$ -trans to old starts

New final state with  $\varepsilon$ -trans from old finals

 $(a|\varepsilon)(c|d|b)^*$ concat (c|d|b)\*a c|d|b d|b

### Alternation:

New start state with  $\varepsilon$ -trans to old starts

New final state with  $\varepsilon$ -trans from old finals

a

 $(a|\varepsilon)(c|d|b)^*$ concat (c|d|b)\* c|d|b d|b

### Alternation:

 $\mathcal{E}$ 

 $S_{M}$ 

New start state with  $\varepsilon$ -trans to old starts

New final state with  $\varepsilon$ -trans from old finals

 $(a|\varepsilon)(c|d|b)^*$ concat (c|d|b)\* a c|d|b d|b Ν M

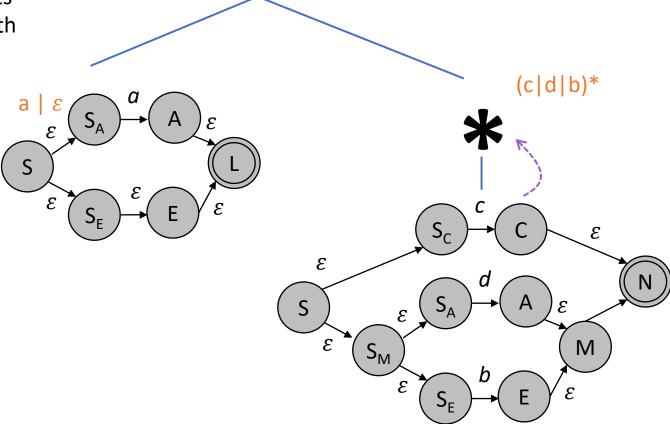
concat

### Alternation:

New start state with  $\varepsilon$ -trans to old starts

New final state with  $\varepsilon$ -trans from old finals

 $(a|\varepsilon)(c|d|b)^*$ 



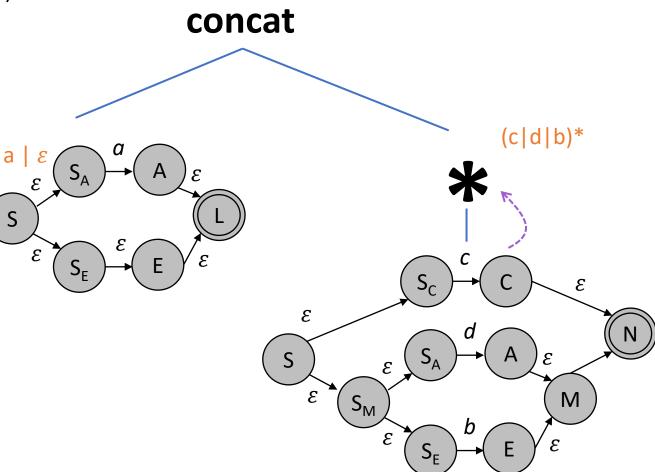
 $(a|\varepsilon)(c|d|b)^*$ 

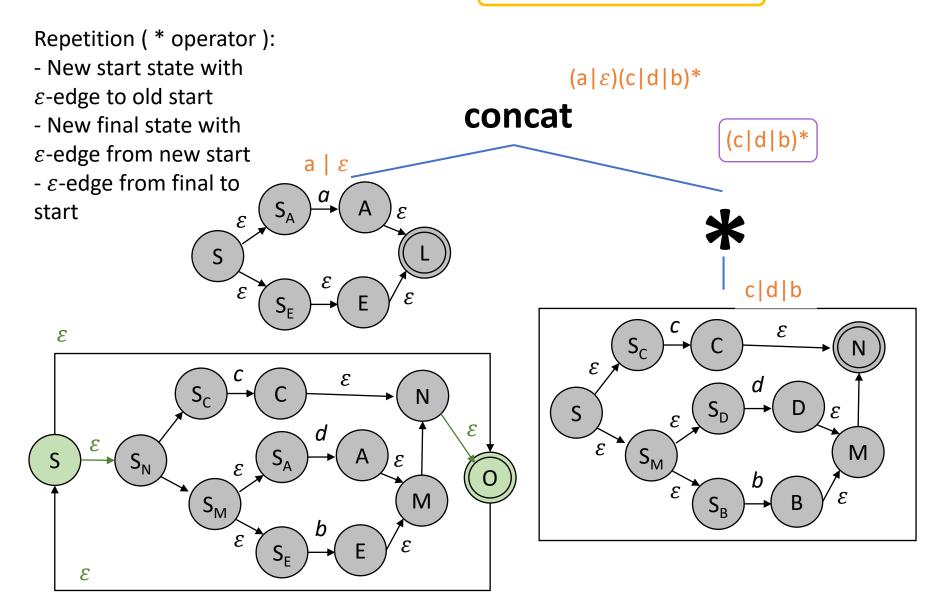
Repetition ( \* operator ):

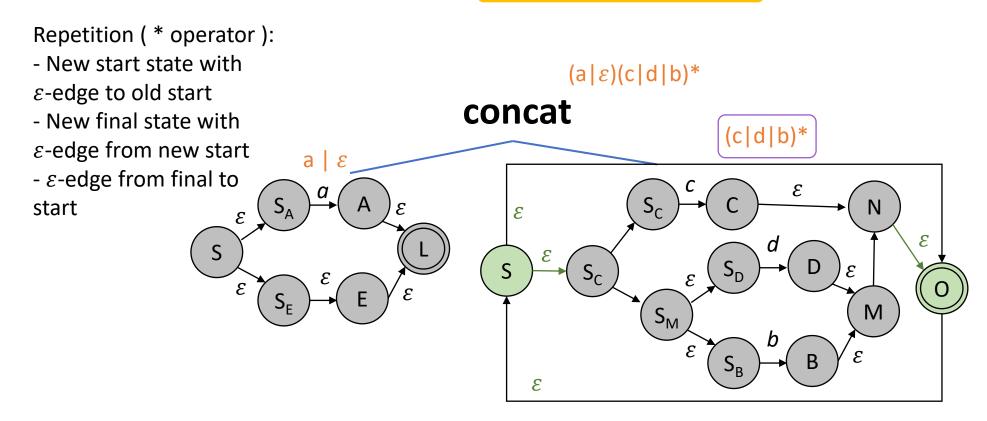
- New start state with  $\varepsilon$ -edge to old start

- New final state with  $\varepsilon$ -edge from new start

-  $\varepsilon$ -edge from final to start







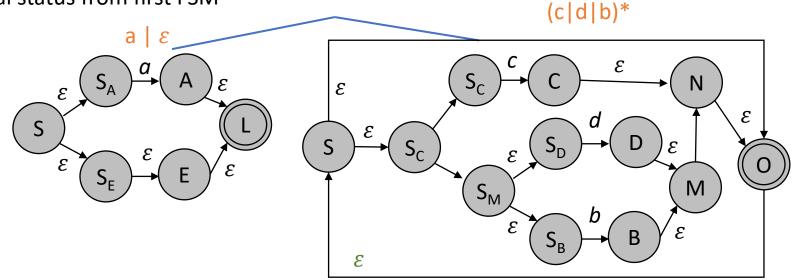
#### Concatenation:

Add  $\varepsilon$ -edge from first FSM's final state to second FSM's start state

Remove final status from first FSM

 $(a|\varepsilon)(c|d|b)^*$ 

concat



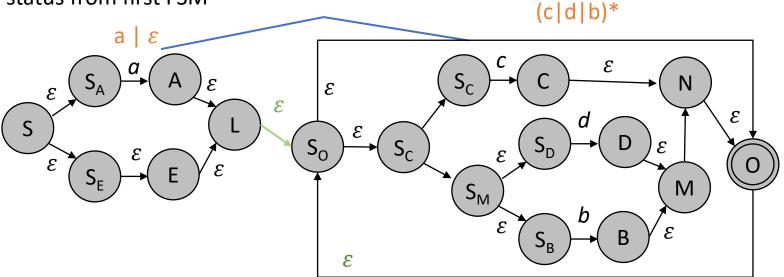
#### Concatenation:

Add  $\varepsilon$ -edge from first FSM's final state to second FSM's start state

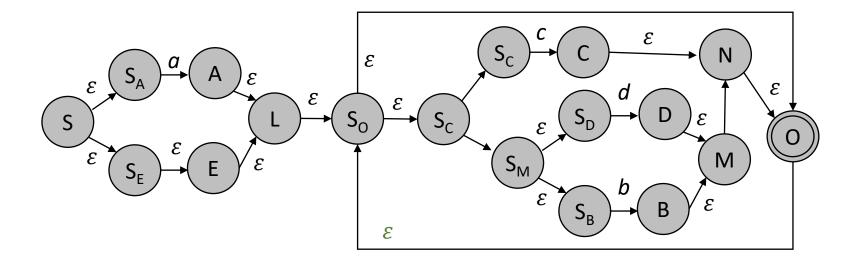
Remove final status from first FSM

 $(a|\varepsilon)(c|d|b)^*$ 

concat



 $(a|\varepsilon)(c|d|b)^*$ 



### Thompson's Construction: Side-Note

Build the RegEx Tree | Replace nodes bottom-up

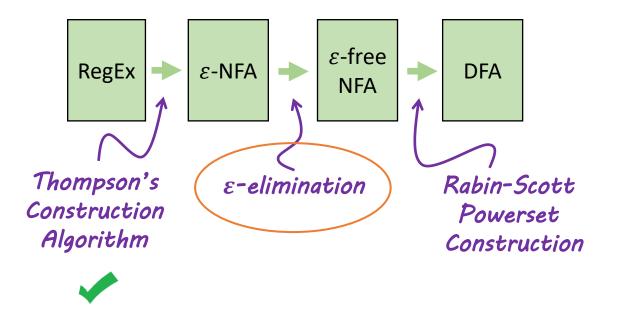
## The FSMs produced by Thompsons Construction are a little bit messy!

- Clearly less efficient than what we would do by hand
- Designed for ease of proofs
- In practice, it's easy to minimize FSMs later



### From RegEx to DFA

Lecture 2 – Implementing Scanners



### Eliminating $\varepsilon$ -transitions

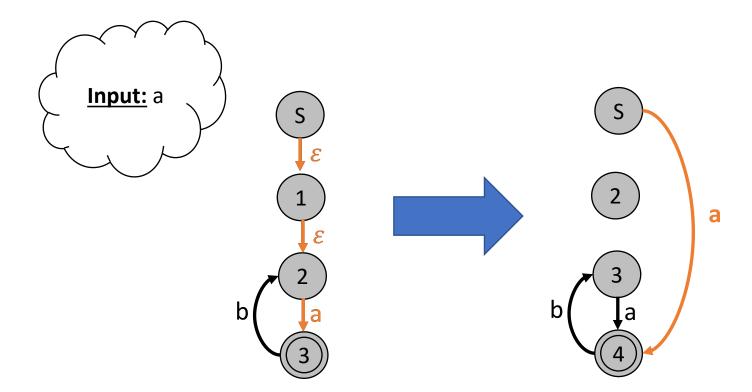
Lecture 2 – Implementing Scanners

### Observation: You never see an epsilon in the input

• Consuming a character means taking a "chain" of zero-or-more  $\varepsilon$ -edges then a real character edge

### **Algorithm Intuition:** cut out the middleman

Replace all "chains" with a direct real-character edge



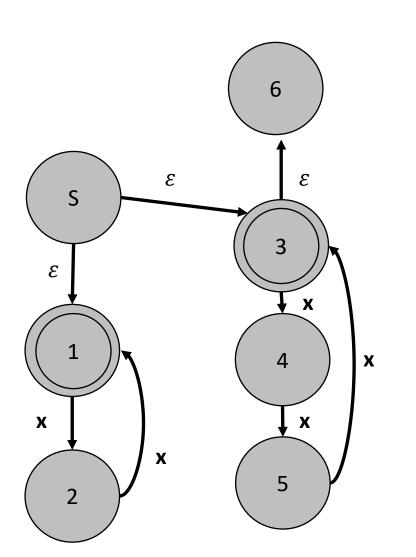
### Eliminating *E*-transitions

Lecture 2 – Implementing Scanners

- Compute  $\varepsilon$ -close(s), the set of states reachable via 0 or more  $\varepsilon$ -edges from s
- Copy all states from N to an  $\varepsilon$ -free version, N'
- Put s in F' if  $\varepsilon$ -close(s) contains a state in F
- Put s,c  $\rightarrow$  t in  $\delta$ ' if there is a c-edge to t in  $\varepsilon$ -close(s)

## Example, Step I Eliminating $\varepsilon$ -Transitions

Let  $\varepsilon$ -close(s) be the set of states reachable via 0 or more  $\varepsilon$ -edges



$$\varepsilon$$
-close(S) = {S,1,3,6}

$$\varepsilon$$
-close(3) = {3, 6}

$$\varepsilon$$
-close(1) = {1}

$$\varepsilon$$
-close(2) = {2}

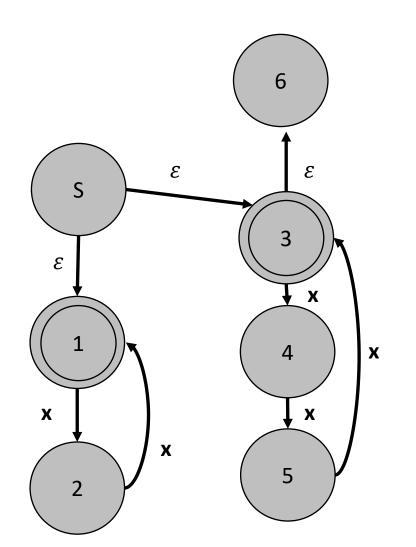
$$\varepsilon$$
-close(4) = {4}

$$\varepsilon$$
-close(5) = {5}

$$\varepsilon$$
-close(6) = {6}

## Example, Step II Eliminating $\varepsilon$ -Transitions

### Copy all states from N to N'



$$\varepsilon$$
-close(S) = {S,1,3,6}

$$\varepsilon$$
-close(3) = {3, 6}

$$\varepsilon$$
-close(1) = {1}

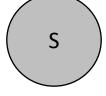
$$\varepsilon$$
-close(2) = {2}

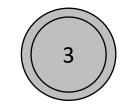
$$\varepsilon$$
-close(4) = {4}

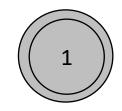
$$\varepsilon$$
-close(5) = {5}

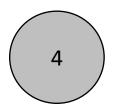
$$\varepsilon$$
-close(6) = {6}

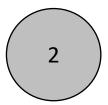


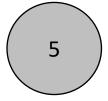






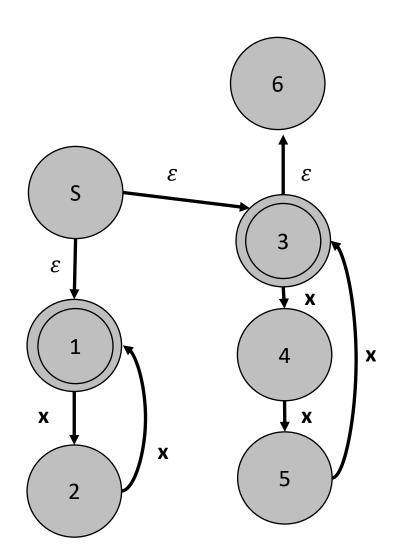






## Example, Step III Eliminating $\varepsilon$ -Transitions

### Put s in F' if $\varepsilon$ -close(s) contains a state in F



$$\varepsilon$$
-close(S) = {S,1,3,6}

$$\varepsilon$$
-close(3) = {3, 6}

$$\varepsilon$$
-close(1) = {1}

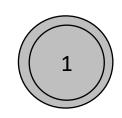
$$\varepsilon$$
-close(2) = {2}

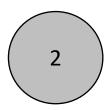
$$\varepsilon$$
-close(4) = {4}

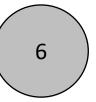
$$\varepsilon$$
-close(5) = {5}

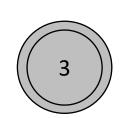
$$\varepsilon$$
-close(6) = {6}

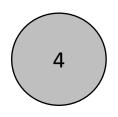


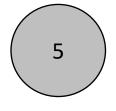






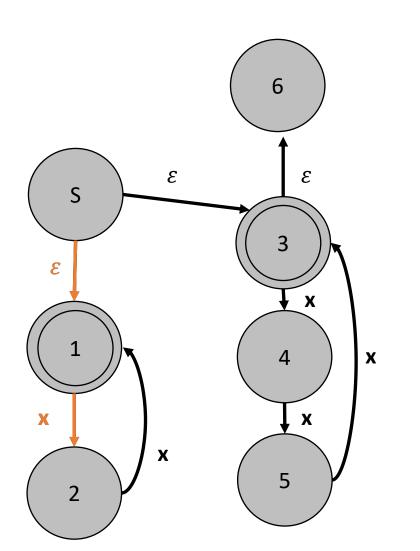






## Example, Step IV Eliminating ε-Transitions

Put s,c  $\rightarrow$  t in  $\delta'$  if there is a c-edge to t in  $\varepsilon$ -close(s)



$$\varepsilon$$
-close(S) = {S,13,6}

$$\varepsilon$$
-close(3) = {3, 6}

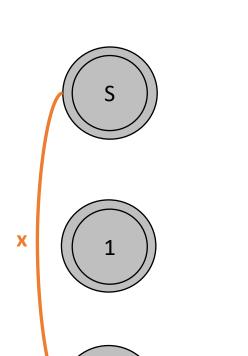
$$\varepsilon$$
-close(1) = {1}

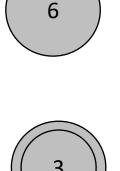
$$\varepsilon$$
-close(2) = {2}

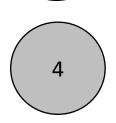
$$\varepsilon$$
-close(4) = {4}

$$\varepsilon$$
-close(5) = {5}

$$\varepsilon$$
-close(6) = {6}



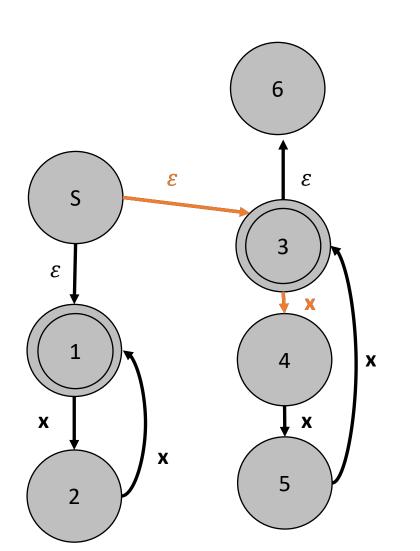






## Example, Step IV Eliminating $\varepsilon$ -Transitions

Put s,c  $\rightarrow$  t in  $\delta'$  if there is a c-edge to t in  $\varepsilon$ -close(s)



$$\varepsilon$$
-close(S) = {S,1,36}

$$\varepsilon$$
-close(3) = {3, 6}

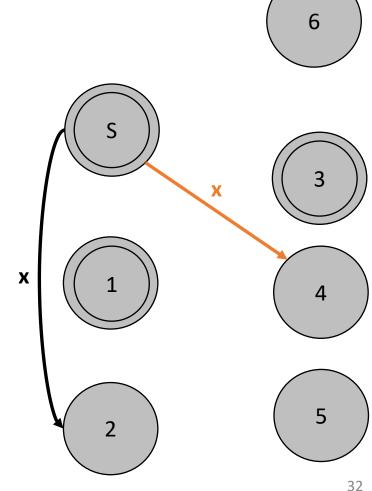
$$\varepsilon$$
-close(1) = {1}

$$\varepsilon$$
-close(2) = {2}

$$\varepsilon$$
-close(4) = {4}

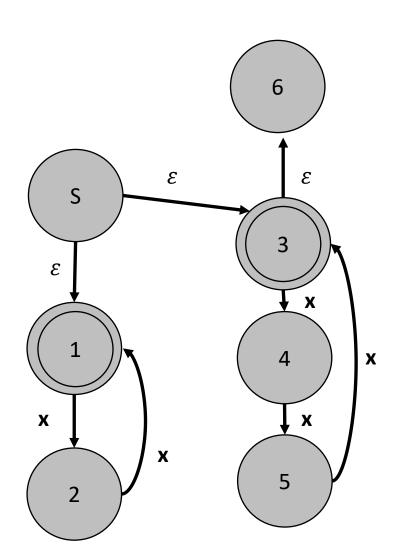
$$\varepsilon$$
-close(5) = {5}

$$\varepsilon$$
-close(6) = {6}



## Example, Step IV Eliminating ε-Transitions

Put s,c  $\rightarrow$  t in  $\delta'$  if there is a c-edge to t in  $\varepsilon$ -close(s)



$$\varepsilon$$
-close(S) = {S,1,3(6)}

$$\varepsilon$$
-close(3) = {3, 6}

$$\varepsilon$$
-close(1) = {1}

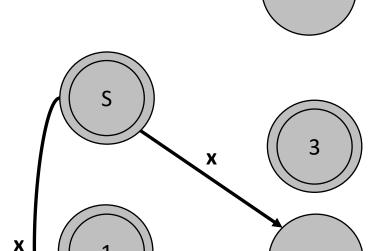
$$\varepsilon$$
-close(2) = {2}

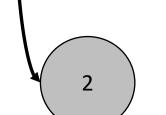
$$\varepsilon$$
-close(4) = {4}

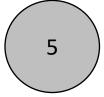
$$\varepsilon$$
-close(5) = {5}

$$\varepsilon$$
-close(6) = {6}







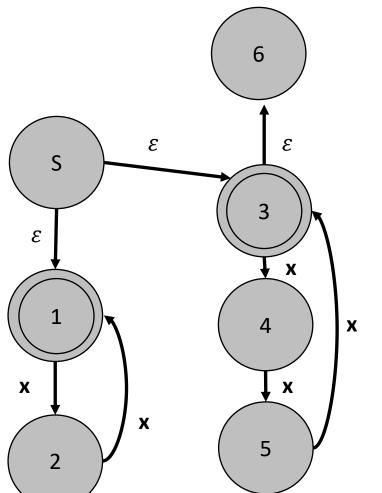


6

## Example, Step IV Eliminating $\epsilon$ -Transitions

### Put s,c $\rightarrow$ t in $\delta'$ if there is a c-edge to t in $\varepsilon$ -close(s)





$$\varepsilon$$
-close(S) = {S,1,3,6}

$$\varepsilon$$
-close(3) = {3, 6}

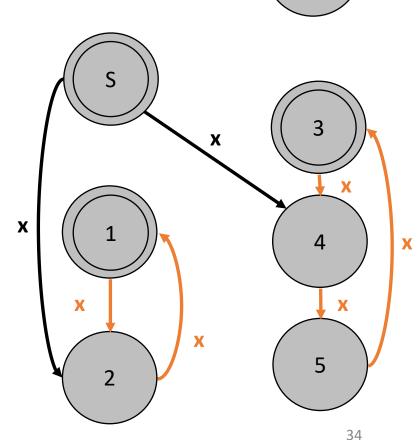
$$\varepsilon$$
-close(1) = {1}

$$\varepsilon$$
-close(2) = {2}

$$\varepsilon$$
-close(4) = {4}

$$\varepsilon$$
-close(5) = {5}

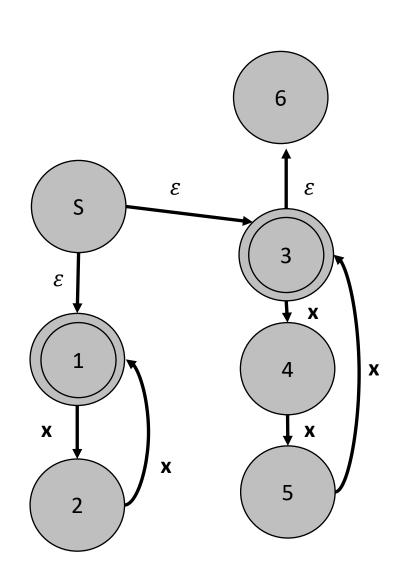
$$\varepsilon$$
-close(6) = {6}

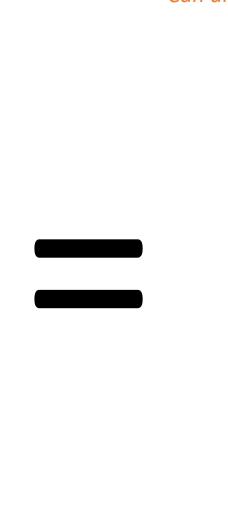


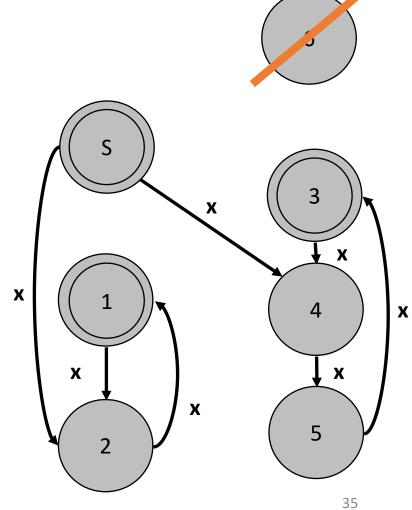
6

## Example, Done! Eliminating ε-Transitions

Can also remove unreachable "useless" state

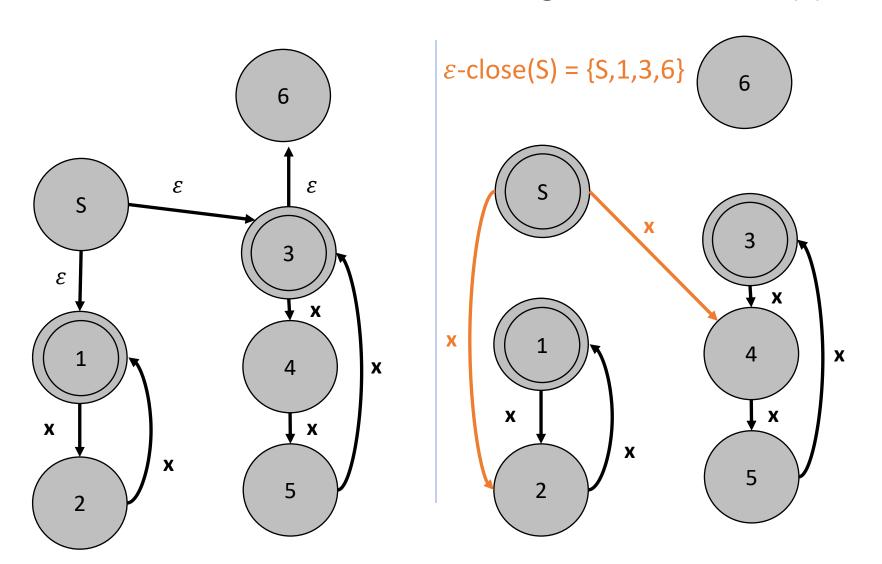






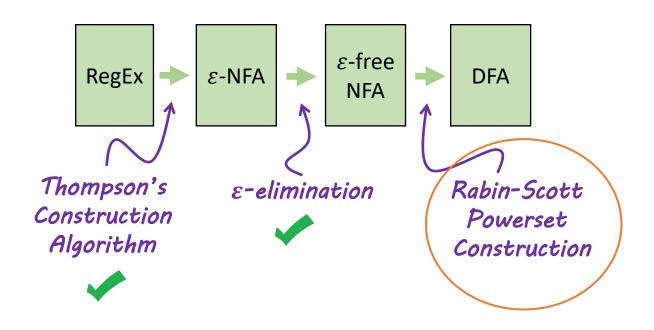
## Example, Step IV Eliminating $\varepsilon$ -Transitions

Put s,c  $\rightarrow$  t in  $\delta'$  if there is a c-edge to t in  $\varepsilon$ -close(s)



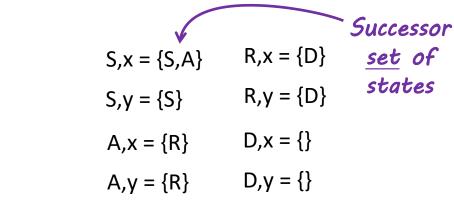
### From RegEx to DFA

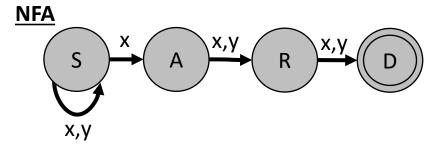
Lecture 2 – Implementing Scanners

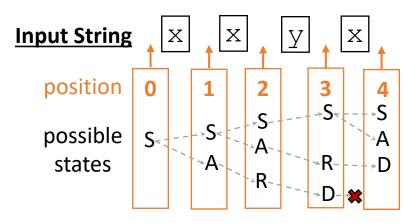


### Recall: NFA Matching Procedure

- NFA can "choose" which transition to take
  - Always moves to states that leads to acceptance (if possible)
- Simulate <u>set</u> of states the NFA could be in
  - If any state in the ending set is final, string accepted







$$S,x = \{S,A\}$$

$$A,x = \{R\}$$

$$R,x = \{D\}$$

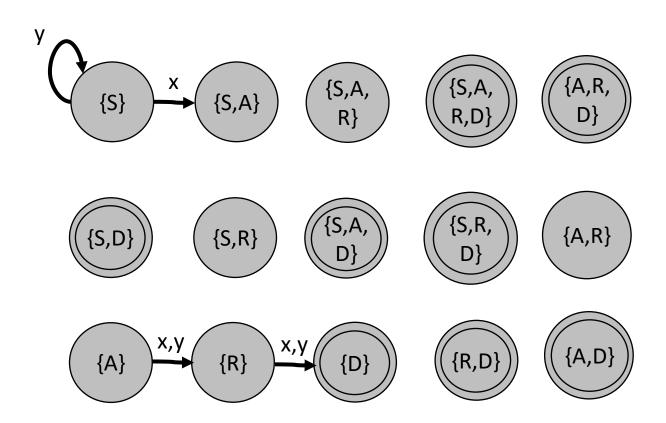
$$D, x = \{\}$$

$$S,y = \{S\}$$

$$S,y = \{S\}$$
  $A,y = \{R\}$ 

$$R,y = \{D\}$$

$$D,y = \{\}$$



$$S, x = \{S,A\} \qquad A, x = \{R\} \qquad R, x = \{D\} \qquad D, x = \{\}$$

$$S, y = \{S\} \qquad A, y = \{R\} \qquad R, y = \{D\} \qquad D, y = \{\}$$

$$\{S,A\}, x = S, x \cup A, x$$

$$= \{S,A\} \cup \{R\}$$

$$= \{S,A,R\}$$

$$\{S,A,R\}$$

$$S, x = \{S,A\} \qquad A, x = \{R\} \qquad R, x = \{D\} \qquad D, x = \{\}$$

$$S, y = \{S\} \qquad A, y = \{R\} \qquad R, y = \{D\} \qquad D, y = \{\}$$

$$\{S,A\}, y = S, y \cup A, y = \{S\} \cup \{R\} = \{S,R\}$$

$$= \{S,R\}$$

$$\{S,A, X \in \{S,A, X \in \{S,A, X \in \{S,A, R,D\}\}\}$$

$$\{S,A\} \qquad X \in \{S,A, R,D\}$$

$$\{S,A\} \qquad X \in \{S,A, R,D\}$$

$$\{S,A\} \qquad X \in \{S,A, R,D\}$$

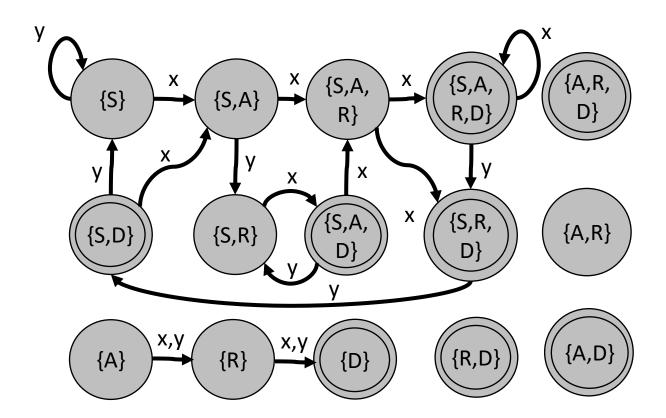
$$\{S,R\} \qquad \{S,A, R,D\}$$

$$\{S,R\} \qquad \{S,R\} \qquad \{S,R,D\}$$

$$\{A,R\} \qquad \{A,R\}$$

$$S,x = \{S,A\}$$
  $A,x = \{R\}$   $R,x = \{D\}$   $D,x = \{\}$ 

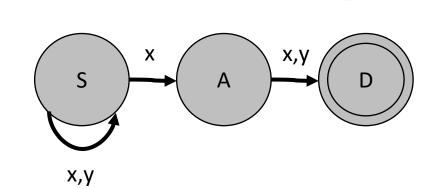
$$S,y = \{S\}$$
  $A,y = \{R\}$   $R,y = \{D\}$   $D,y = \{\}$ 



## Exponential State Count Rabin-Scott Powerset Construction

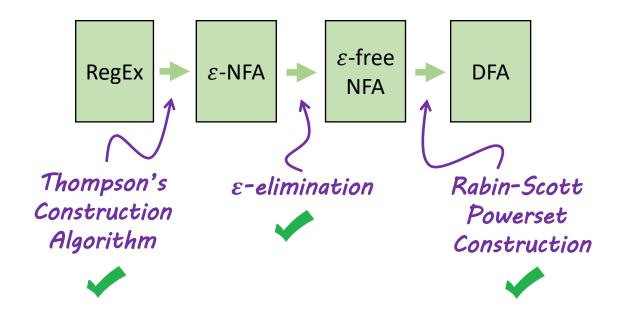
- How may states might the DFA have?
  - · 2|Q|
- Why 2 |Q|?

<u>S</u>	<u>A</u>	<u>D</u>	
0	0	0	{}
0	0	1	{D}
0	1	0	{A}
1	0	0	{S}
0	1	1	{A,D}
1	1	0	{S,A}
1	0	1	{S,D}
1	1	1	{S,A,D}



### From RegEx to DFA

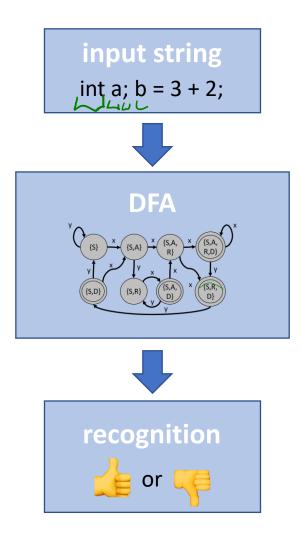
Lecture 2 – Implementing Scanners



## DONE! ... or are we?

## DFA # Tokenizer Limitations

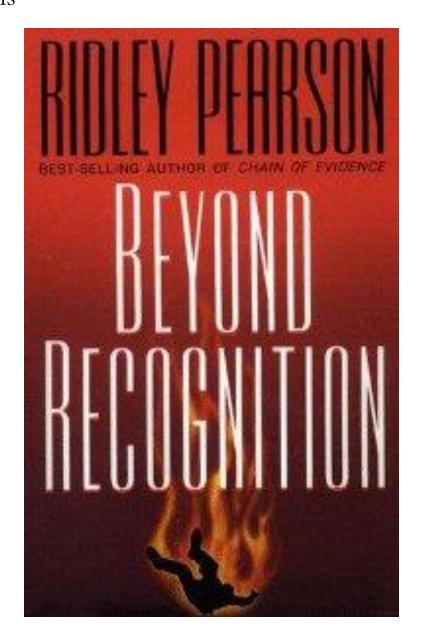
- Finite automata only check for language membership of a string (recognition)
- The Scanner needs to
  - Break the input into many different tokens
  - Know what characters comprise the token



### DFA # Tokenizer Limitations

- Finite automata only check for language membership of a string (recognition)
- The Scanner needs to
  - Break the input into many different tokens
  - Know what characters comprise the token

We need to go... beyond recognition



## Next Time Lecture 3 Preview

