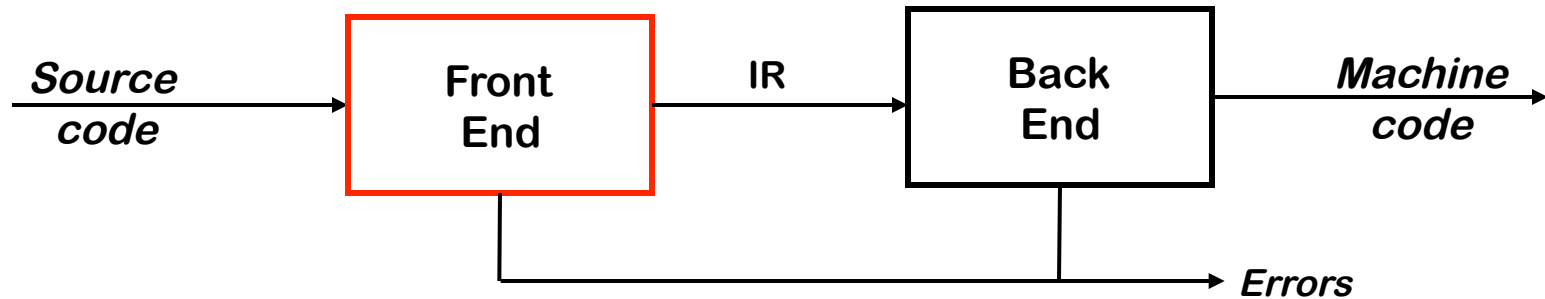


# *CS415 Compilers*

## *Lexical Analysis*

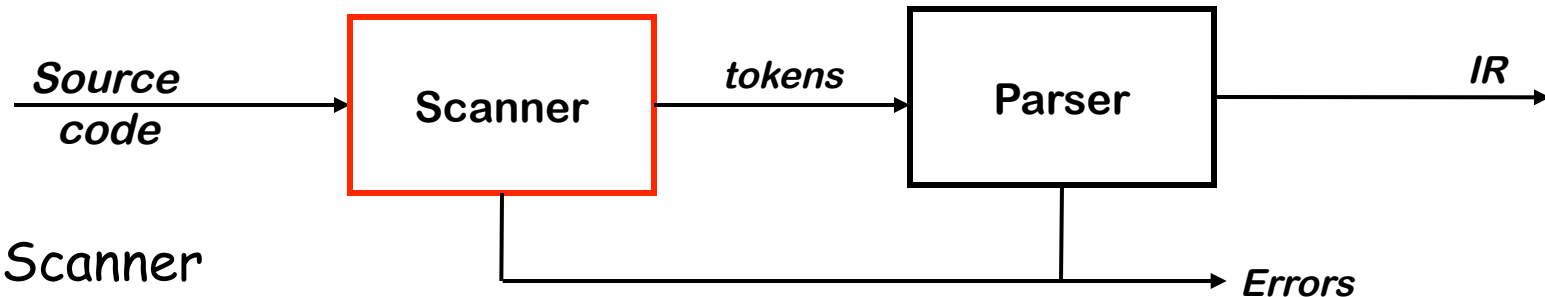
These slides are based on slides copyrighted by  
Keith Cooper, Ken Kennedy & Linda Torczon at Rice  
University



The purpose of the front end is to deal with the input language

- Perform a membership test:  $\text{code} \in \text{source language?}$
- Is the program well-formed (semantically) ?
- Build an IR version of the code for the rest of the compiler

*The front end is not monolithic*

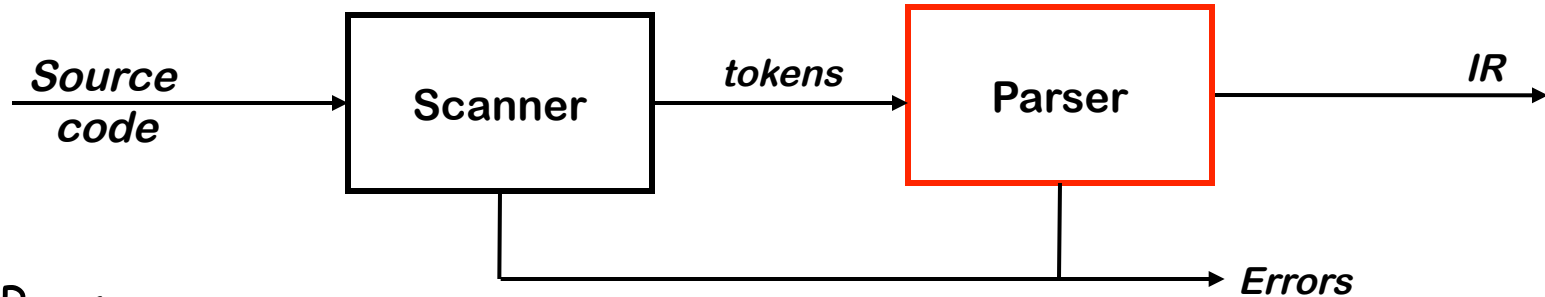


### Scanner

- Maps stream of characters into words
  - Basic unit of syntax
  - $x = x + y ;$  becomes  
 $\langle \text{id}, x \rangle \langle \text{eq}, = \rangle \langle \text{id}, x \rangle \langle \text{pl}, + \rangle \langle \text{id}, y \rangle \langle \text{sc}, ; \rangle$

Speed is an issue in scanning  
 ⇒ use a specialized recognizer

- Characters that form a word are its *lexeme*
- Its *part of speech* (or *syntactic category*) is called its *token type*
- Scanner discards white space & (often) comments



### Parser

- Checks stream of classified words (*parts of speech*) for grammatical correctness
- Determines if code is syntactically well-formed
- Guides checking at deeper levels than syntax
- Builds an IR representation of the code

*We'll get to parsing in the next lectures*

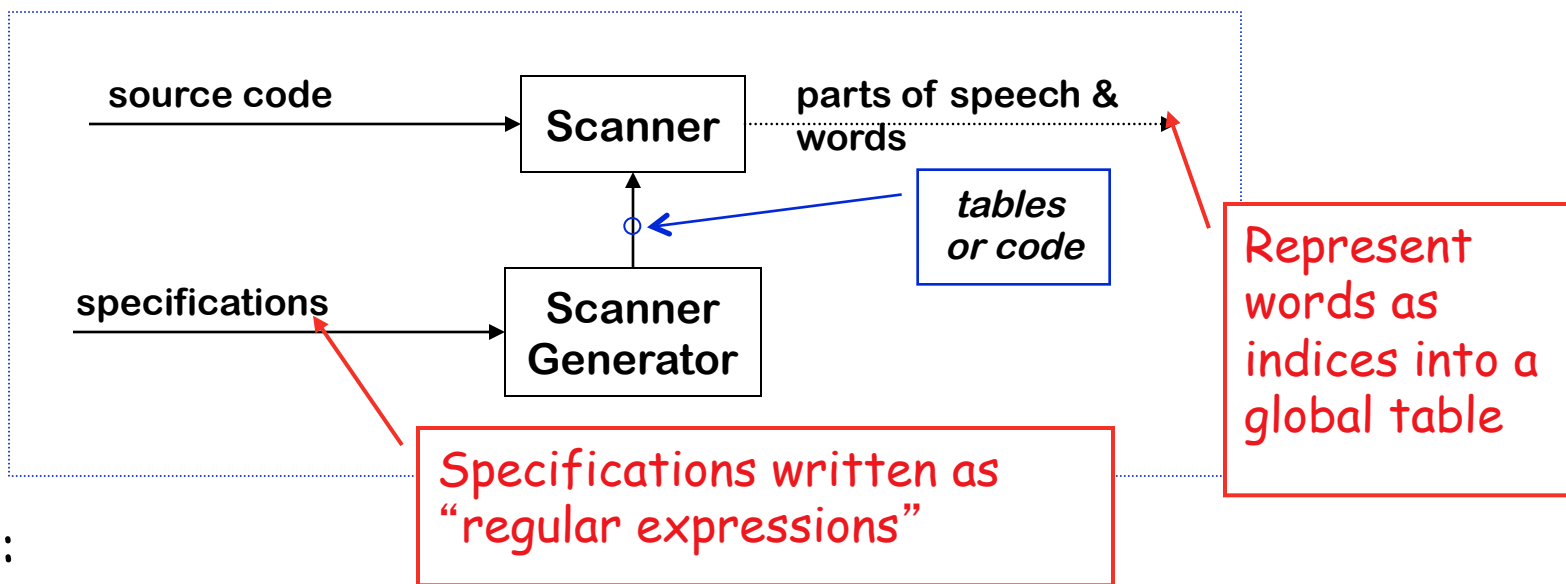
- Language syntax is specified with *parts of speech*, not *words*
- Syntax checking matches *parts of speech* against a grammar

1.  $goal \rightarrow expr$   
2.  $expr \rightarrow expr\ op\ term$   
3.       |  $term$   
4.  $term \rightarrow \underline{number}$   
5.       |  $\underline{id}$   
6.  $op \rightarrow +$   
7.       |  $-$

$S = goal$   
 $T = \{ \underline{number}, \underline{id}, +, - \}$   
 $N = \{ goal, expr, term, op \}$   
 $P = \{ 1, 2, 3, 4, 5, 6, 7 \}$

Why study lexical analysis?

- We want to avoid writing scanners by hand



Goals:

- To simplify specification & implementation of scanners
- To understand the underlying techniques and technologies

Lexical patterns form a *regular language*

*\*\*\* any finite language is regular \*\*\**

*Regular expressions* (REs) describe regular languages

Ever type  
"rm \*.o a.out" ?

Regular Expression (over alphabet  $\Sigma$ )

- $\varepsilon$  is a RE denoting the set  $\{\varepsilon\}$
- If  $a$  is in  $\Sigma$ , then  $a$  is a RE denoting  $\{a\}$
- If  $x$  and  $y$  are REs denoting  $L(x)$  and  $L(y)$  then
  - $x | y$  is an RE denoting  $L(x) \cup L(y)$
  - $xy$  is an RE denoting  $L(x)L(y)$
  - $x^*$  is an RE denoting  $L(x)^*$

Precedence is  
closure, then  
concatenation,  
then alternation

Operation	Definition
<i>Union of L and M</i> Written $L \cup M$	$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$
<i>Concatenation of L and M</i> Written $LM$	$LM = \{st \mid s \in L \text{ and } t \in M\}$
<i>Kleene closure of L</i> Written $L^*$	$L^* = \bigcup_{0 \leq i < \infty} L^i$
<i>Positive Closure of L</i> Written $L^+$	$L^+ = \bigcup_{1 \leq i < \infty} L^i$

*These definitions should be well known*

## Identifiers:

*Letter* → (a|b|c| ... |z|A|B|C| ... |Z)

*Digit* → (0|1|2| ... |9)

*Identifier* → *Letter* (*Letter* | *Digit*)\*

## Numbers:

*Integer* → (+|-|ε) (0| (1|2|3| ... |9)(*Digit*\*) )

*Decimal* → *Integer* .*Digit*\*

*Real* → ( *Integer* | *Decimal* ) E (+|-|ε) *Digit*\*

*Complex* → ( *Real* , *Real* )

*Numbers can get much more complicated!*

- All strings of 1s and 0s ending in a 1

$(\underline{0} | \underline{1})^* \underline{1}$

- All strings over lowercase letters where the vowels (a,e,i,o, & u) occur exactly once, in ascending order

*Cons*  $\rightarrow$   $(\underline{b} | \underline{c} | \underline{d} | \underline{f} | \underline{g} | \underline{h} | \underline{j} | \underline{k} | \underline{l} | \underline{m} | \underline{n} | \underline{p} | \underline{q} | \underline{r} | \underline{s} | \underline{t} | \underline{v} | \underline{w} | \underline{x} | \underline{y} | \underline{z})$

*Cons*<sup>\*</sup>  $\underline{a}$  *Cons*<sup>\*</sup>  $\underline{e}$  *Cons*<sup>\*</sup>  $\underline{i}$  *Cons*<sup>\*</sup>  $\underline{o}$  *Cons*<sup>\*</sup>  $\underline{u}$  *Cons*<sup>\*</sup>

- All strings of 1s and 0s that do not contain three 1s in a row:

- All strings of 1s and 0s ending in a 1

$(\underline{0} | \underline{1})^* \underline{1}$

- All strings over lowercase letters where the vowels (a,e,i,o, & u) occur exactly once, in ascending order

$Cons \rightarrow (\underline{b} | \underline{c} | \underline{d} | \underline{f} | \underline{g} | \underline{h} | \underline{j} | \underline{k} | \underline{l} | \underline{m} | \underline{n} | \underline{p} | \underline{q} | \underline{r} | \underline{s} | \underline{t} | \underline{v} | \underline{w} | \underline{x} | \underline{y} | \underline{z})$

$Cons^* \underline{a} Cons^* \underline{e} Cons^* \underline{i} Cons^* \underline{o} Cons^* \underline{u} Cons^*$

- All strings of 1s and 0s that do not contain three 1s in a row:

$(\underline{0}^* (\varepsilon | \underline{10} | \underline{110}) \underline{0}^*)^* (\varepsilon | \underline{1} | \underline{11})$

*Regular expressions can be used to specify the words to be translated to parts of speech by a lexical analyzer*

Using results from automata theory and theory of algorithms, we can automatically build recognizers from regular expressions

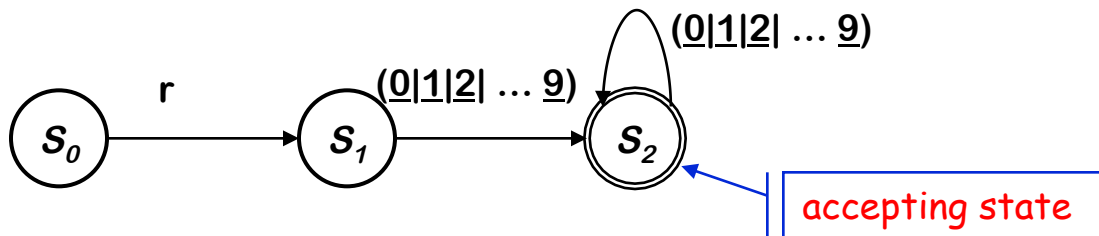
⇒ We study REs and associated theory to automate scanner construction !

Consider the problem of recognizing ILOC register names

$Register \rightarrow r (0|1|2| \dots | 9) (0|1|2| \dots | 9)^*$

- Allows registers of arbitrary number
- Requires at least one digit

RE corresponds to a recognizer (or DFA)

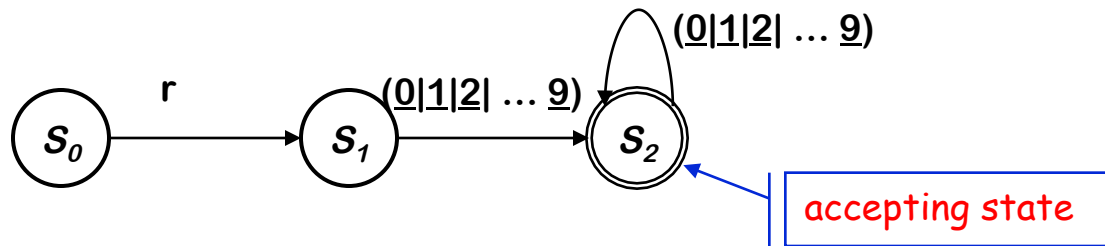


Recognizer for *Register*

*Transitions on other inputs go to an error state,  $s_e$*

## DFA operation

- Start in state  $S_0$  & take transitions on each input character
- DFA accepts a word  $\underline{x}$  iff  $\underline{x}$  leaves it in a final state ( $S_2$ )



So,

Recognizer for *Register*

- r17 takes it through  $s_0, s_1, s_2$  and accepts
- r takes it through  $s_0, s_1$  and fails
- a takes it straight to error state  $s_e$  (not shown here)

To be useful, recognizer must turn into code

```

Char ← next character
State ← s0
while (Char ≠ EOF)
  State ← δ(State,Char)
  Char ← next character
if (State is a final state)
  then report success
  else report failure
  
```

*Skeleton recognizer*

$\delta$	$r$	0,1,2,3,4, 5,6,7,8,9	All others
$s_0$	$s_1$	$s_e$	$s_e$
$s_1$	$s_e$	$s_2$	$s_e$
$s_2$	$s_e$	$s_2$	$s_e$
$s_e$	$s_e$	$s_e$	$s_e$

*Table encoding RE*

To be useful, recognizer must turn into code

```

Char ← next character
State ← s0
while (Char ≠ EOF)
  State ← δ(State,Char)
  perform specified action
  Char ← next character
if (State is a final state)
  then report success
  else report failure
  
```

*Skeleton recognizer*

$\delta$	r	0,1,2,3,4, 5,6,7,8,9	All others
s <sub>0</sub>	s <sub>1</sub> <i>start</i>	s <sub>e</sub> <i>error</i>	s <sub>e</sub> <i>error</i>
s <sub>1</sub>	s <sub>e</sub> <i>error</i>	s <sub>2</sub> <i>add</i>	s <sub>e</sub> <i>error</i>
s <sub>2</sub>	s <sub>e</sub> <i>error</i>	s <sub>2</sub> <i>add</i>	s <sub>e</sub> <i>error</i>
s <sub>e</sub>	s <sub>e</sub> <i>error</i>	s <sub>e</sub> <i>error</i>	s <sub>e</sub> <i>error</i>

*Table encoding RE*

r *Digit Digit\** allows arbitrary numbers

- Accepts r00000
- Accepts r99999
- What if we want to limit it to r0 through r31 ?

Write a tighter regular expression

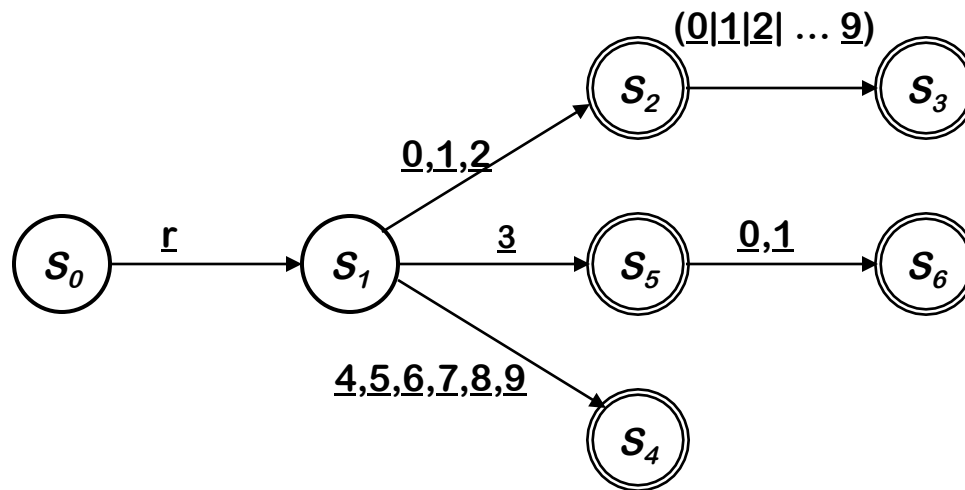
- *Register* → r ( 0|1|2 ) ( *Digit* |  $\epsilon$  ) | ( 4|5|6|7|8|9 ) | ( 3|30|31 ) )
- *Register* → r0|r1|r2| ... |r31|r00|r01|r02| ... |r09

Produces a more complex DFA

- Has more states
- Same cost per transition
- Same basic implementation

The DFA for

*Register*  $\rightarrow \underline{r} ( (\underline{0}|\underline{1}|\underline{2}) (\textit{Digit} | \varepsilon) | (\underline{4}|\underline{5}|\underline{6}|\underline{7}|\underline{8}|\underline{9}) | (\underline{3}|\underline{30}|\underline{31}) )$

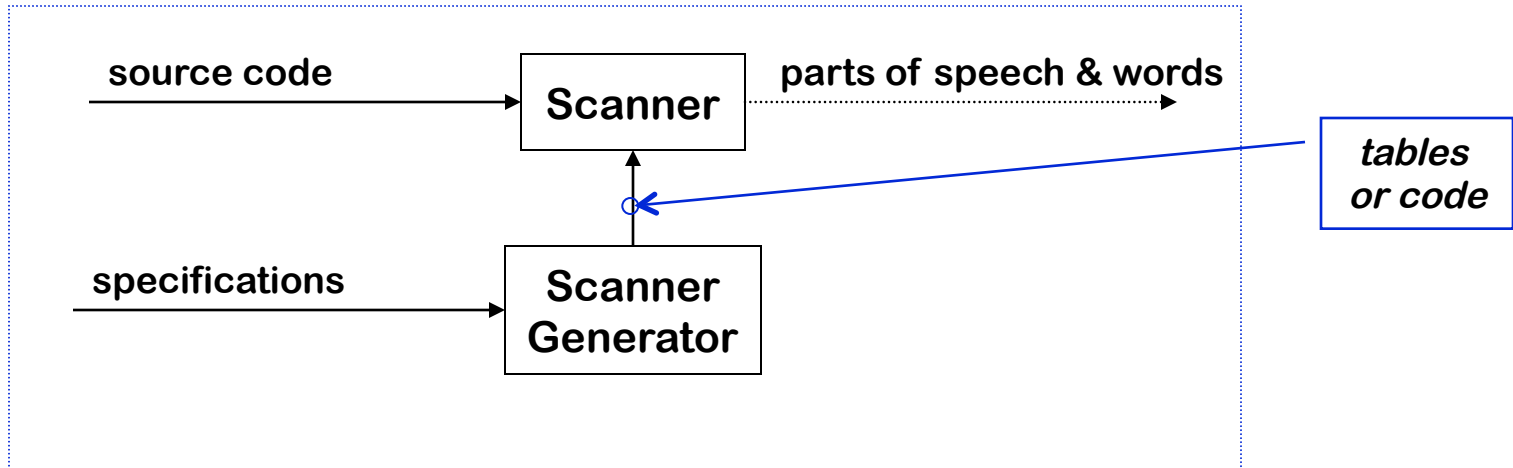


- Accepts a more constrained set of registers
- Same set of actions, more states

$\delta$	r	0,1	2	3	4-9	All others
$s_0$	$s_1$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$
$s_1$	$s_e$	$s_2$	$s_2$	$s_5$	$s_4$	$s_e$
$s_2$	$s_e$	$s_3$	$s_3$	$s_3$	$s_3$	$s_e$
$s_3$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$
$s_4$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$
$s_5$	$s_e$	$s_6$	$s_e$	$s_e$	$s_e$	$s_e$
$s_6$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$
$s_e$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$	$s_e$

Runs in the same skeleton recognizer

*Table encoding RE for the tighter register specification*



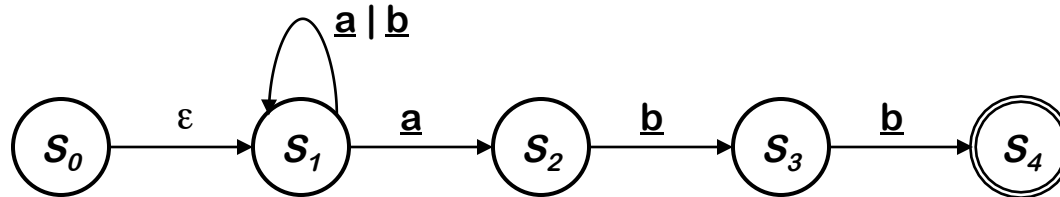
- The scanner is the first stage in the front end
- Specifications can be expressed using regular expressions
- Build tables and code from a DFA

- We will show how to construct a finite state automaton to recognize any RE
- Overview:
  - Direct construction of a **nondeterministic finite automaton (NFA)** to recognize a given RE
    - Requires  $\varepsilon$ -transitions to combine regular subexpressions
  - Construct a **deterministic finite automaton (DFA)** to simulate the NFA
    - Use a set-of-states construction
  - Minimize the number of states
    - Hopcroft state minimization algorithm
  - Generate the scanner code
    - Additional specifications needed for details

Each RE corresponds to a non *deterministic finite automaton* (NFA)

- May be hard to directly construct the right DFA

What about an RE such as  $(\underline{a} \mid \underline{b})^* \underline{a}\underline{b}\underline{b}$  ?



This is a little different

- $S_0$  has a transition on  $\epsilon$
- $S_1$  has two transitions on  $\underline{a}$

This is a *non-deterministic finite automaton* (NFA)

- An NFA accepts a string  $x$  iff  $\exists$  a path through the transition graph from  $s_0$  to a final state such that the edge labels spell  $x$
- Transitions on  $\epsilon$  consume no input
- To “run” the NFA, start in  $s_0$  and *guess* the right transition at each step
  - Always guess correctly
  - If some sequence of correct guesses accepts  $x$  then accept

Why study NFAs?

- They are the key to automating the RE $\rightarrow$ DFA construction
- We can paste together NFAs with  $\epsilon$ -transitions



DFA is a special case of an NFA

- DFA has no  $\epsilon$  transitions
- DFA's transition function is single-valued
- Same rules will work

DFA can be simulated with an NFA

→ *Obviously*

NFA can be simulated with a DFA

*(less obvious)*

- Simulate sets of possible states
- Possible exponential blowup in the state space
- Still, one state per character in the input stream

To convert a specification into code:

- 1 Write down the RE for the input language
- 2 Build a NFA
- 3 Build the DFA that simulates the NFA
- 4 Systematically shrink the DFA
- 5 Turn it into code

Scanner generators

- Lex and Flex work along these lines
- Algorithms are well-known and well-understood
- Key issue is interface to parser *(define all parts of speech)*
- You could build one in a weekend!

RE  $\rightarrow$  NFA (*Thompson's construction*)

- Build an NFA for each term
- Combine them with  $\epsilon$ -moves

NFA  $\rightarrow$  DFA (*subset construction*)

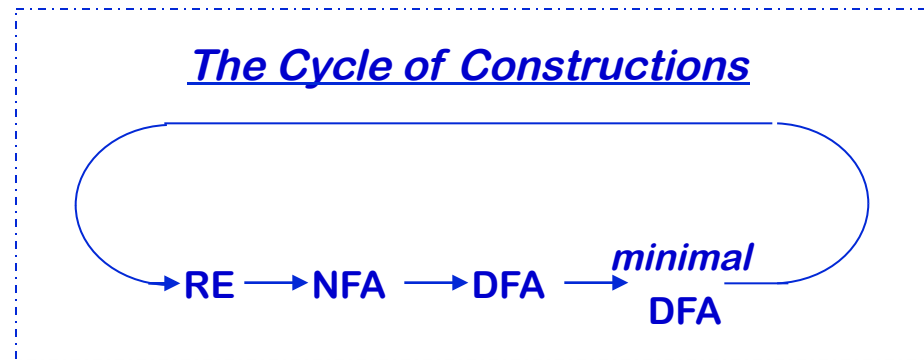
- Build the simulation

DFA  $\rightarrow$  Minimal DFA

- Hopcroft's algorithm

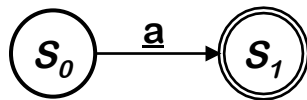
DFA  $\rightarrow$  RE (*Not part of the scanner construction*)

- All pairs, all paths problem
- Take the union of all paths from  $s_0$  to an accepting state



## Key idea

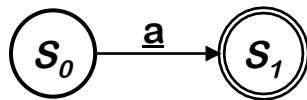
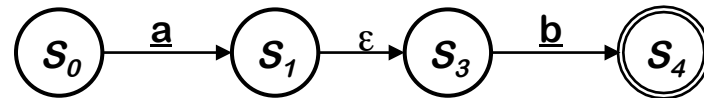
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for ab

Ken Thompson, CACM, 1968

## Key idea

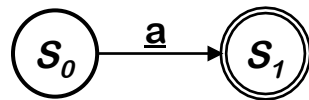
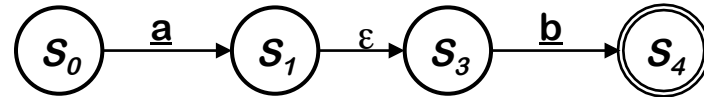
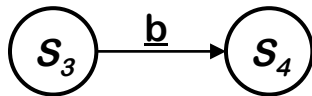
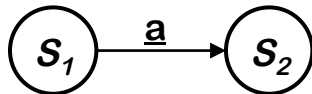
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for ab

Ken Thompson, CACM, 1968

## Key idea

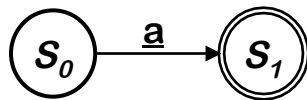
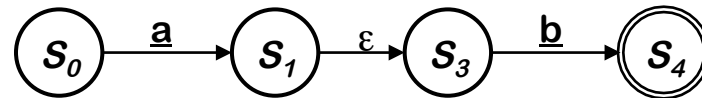
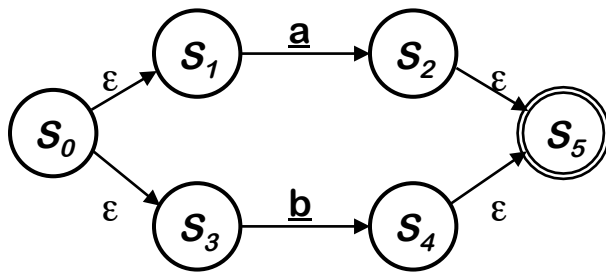
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

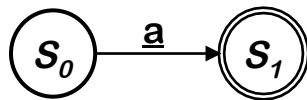
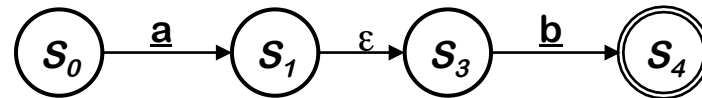
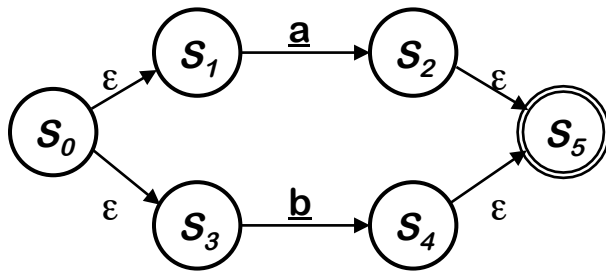
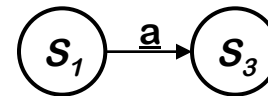
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

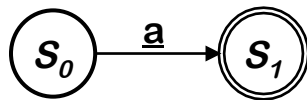
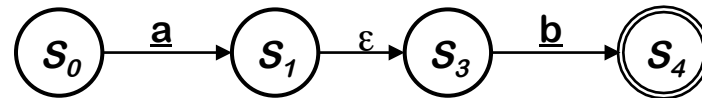
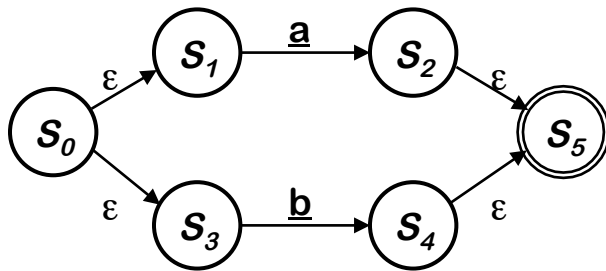
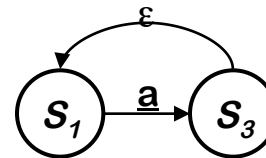
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

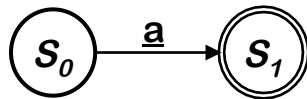
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

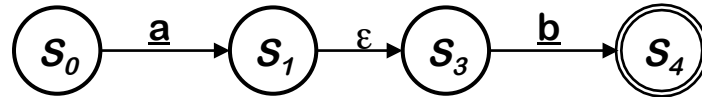
Ken Thompson, CACM, 1968

## Key idea

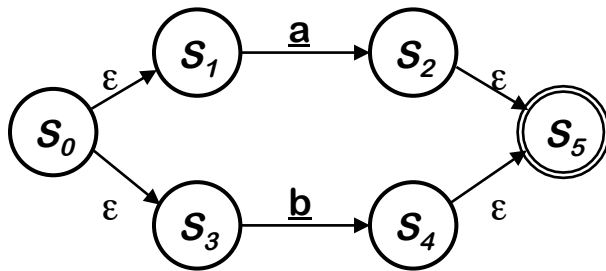
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order



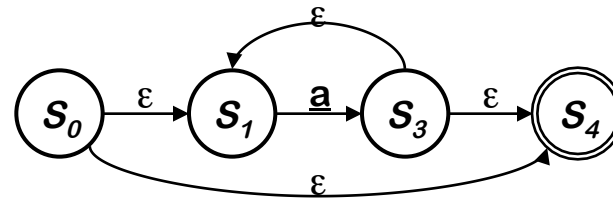
NFA for a



NFA for ab



NFA for a | b

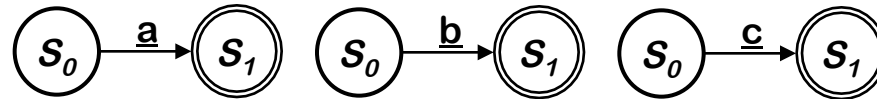


NFA for a\*

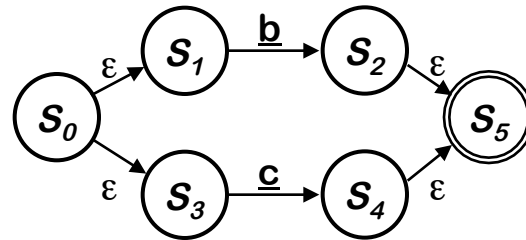
Ken Thompson, CACM, 1968

Let's try  $a(b|c)^*$

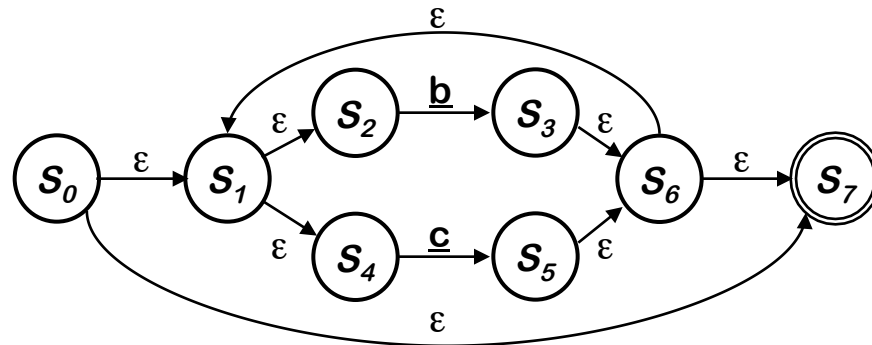
1.  $a, b, \& c$

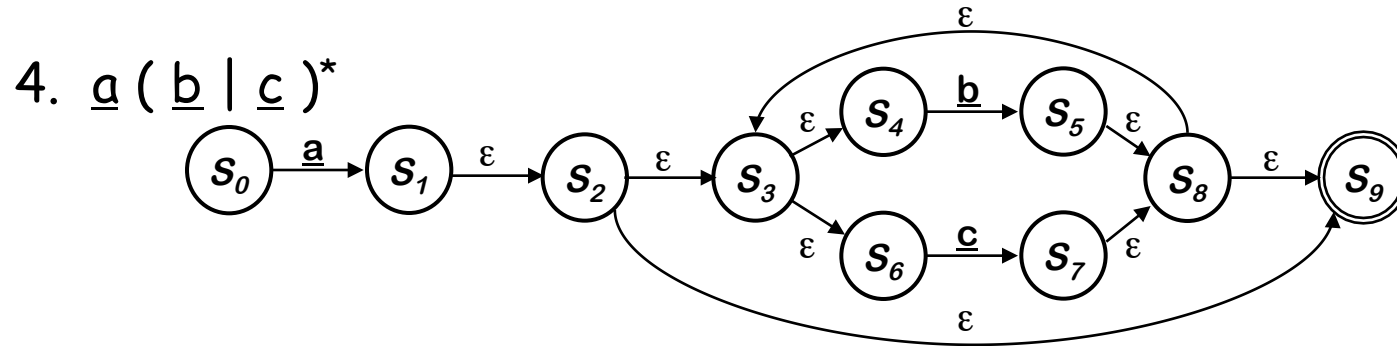


2.  $b|c$

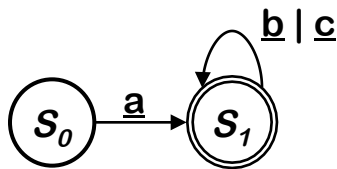


3.  $(b|c)^*$





Of course, a human would design something simpler ...



But, we can automate production of the more complex one ...

Need to build a simulation of the NFA

Two key functions

- $Move(s_i, \underline{a})$  is set of states reachable from  $s_i$  by  $\underline{a}$
- $\varepsilon$ -closure( $s_i$ ) is set of states reachable from  $s_i$  by  $\varepsilon$

The algorithm:

- Start state derived from  $s_0$  of the NFA
- Take its  $\varepsilon$ -closure  $S_0 = \varepsilon$ -closure( $s_0$ )
- Take the image of  $S_0$ ,  $move(S_0, a)$  for each  $a \in \Sigma$ , and take its  $\varepsilon$ -closure
- Iterate until no more states are added

*Sounds more complex than it is...*

**The algorithm:**

```

 $s_0 \leftarrow \varepsilon\text{-closure}(q_0)$ 
add  $s_0$  to  $S$ 
while (  $S$  is still changing )
  for each  $s_i \in S$ 
    for each  $a \in \Sigma$ 
       $s_j \leftarrow \varepsilon\text{-closure}(\text{move}(s_i, a))$ 
      if (  $s_j \notin S$  ) then
        add  $s_j$  to  $S$  as  $s_j$ 
         $T[s_i, a] \leftarrow s_j$ 
      else
         $T[s_i, a] \leftarrow s_j$ 

```

*Let's think about why this works*

**The algorithm halts:**

1.  $S$  contains no duplicates (test before adding)
2.  $2^Q$  is finite
3. while loop adds to  $S$ , but does not remove from  $S$  (*monotone*)

$\Rightarrow$  the loop halts

$S$  contains all the reachable NFA states

*It tries each character in each  $s_i$ .  
It builds every possible NFA configuration.*

$\Rightarrow S$  and  $T$  form the DFA

Example of a *fixed-point* computation

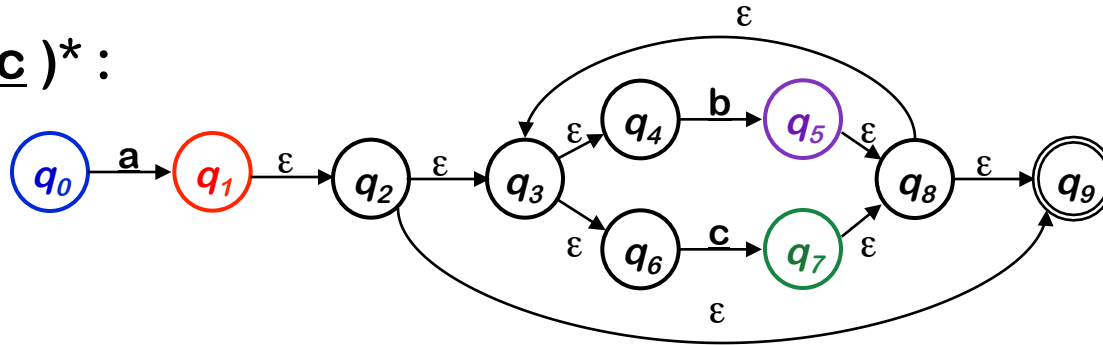
- Monotone construction of some finite set
- Halts when it stops adding to the set
- Proofs of halting & correctness are similar
- These computations arise in many contexts

Other fixed-point computations

- Canonical construction of sets of LR(1) items
  - $\rightarrow$  Quite similar to the subset construction
- Classic data-flow analysis
  - $\rightarrow$  Solving sets of simultaneous set equations
- DFA minimization algorithm (coming up!)

*We will see many more fixed-point computations*

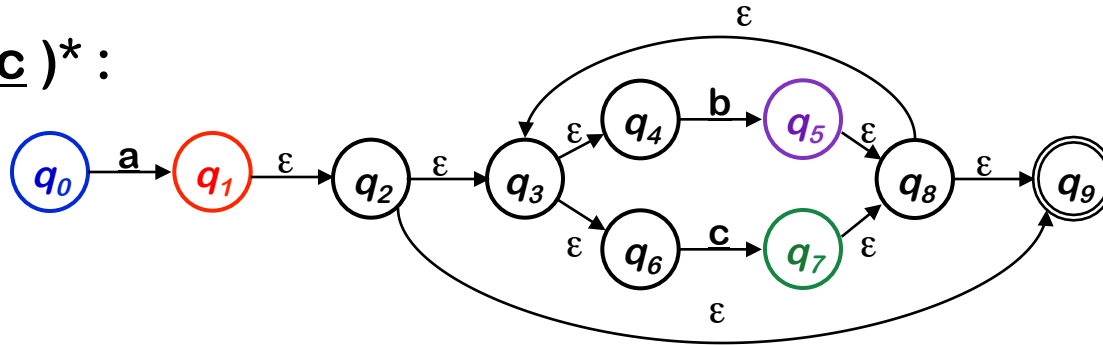
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$			

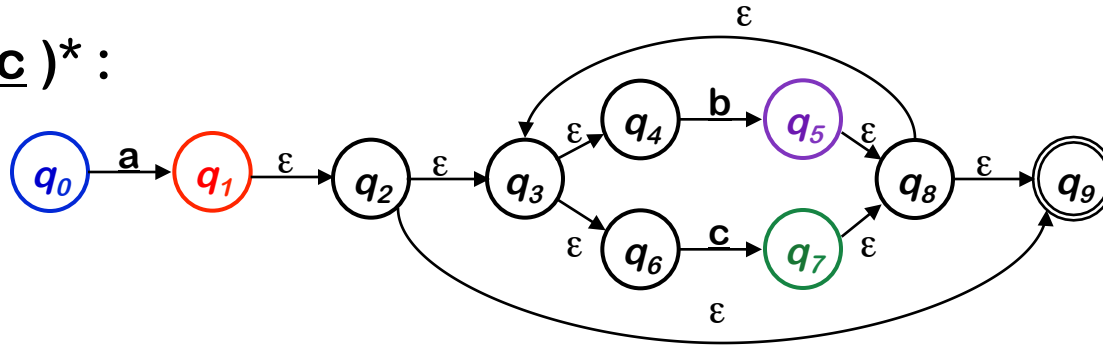
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none

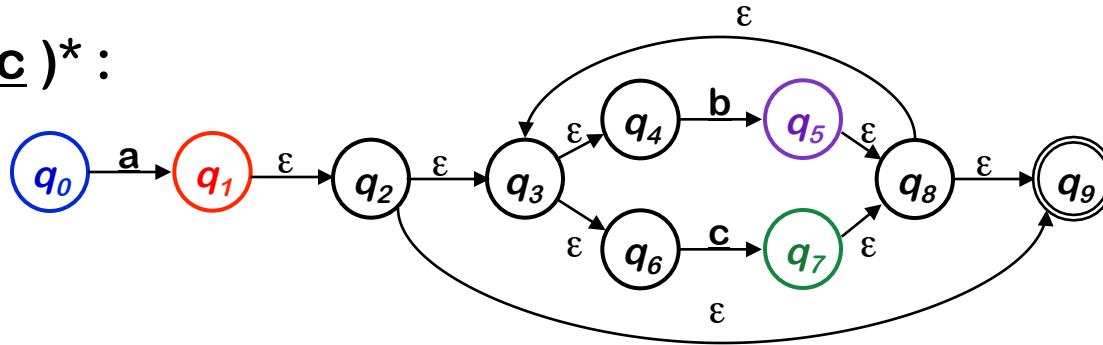
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$			

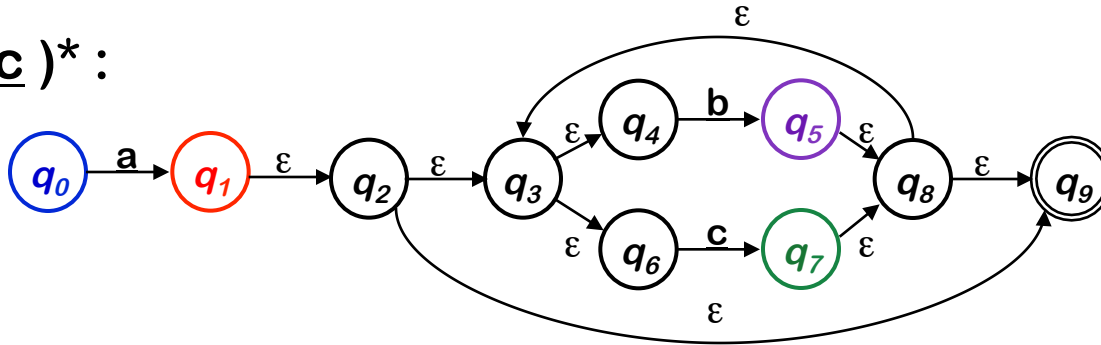
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none		

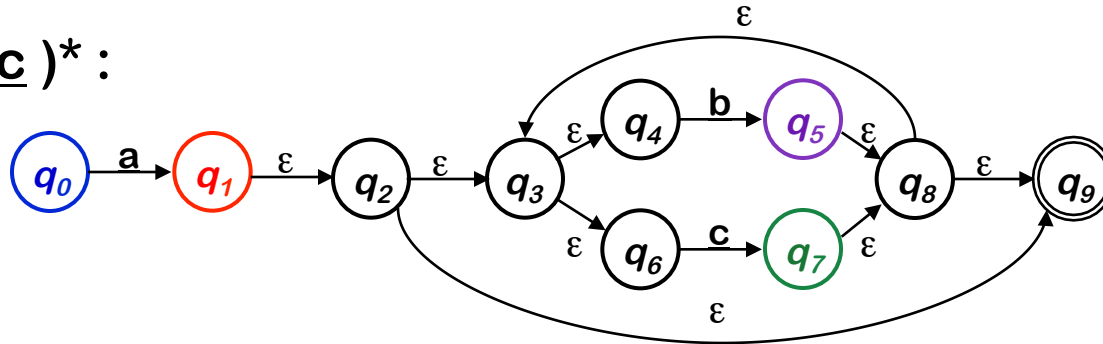
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	...

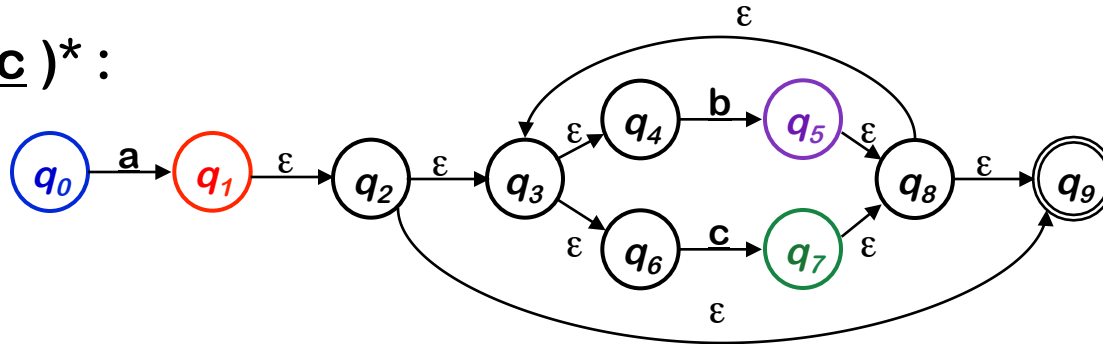
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$

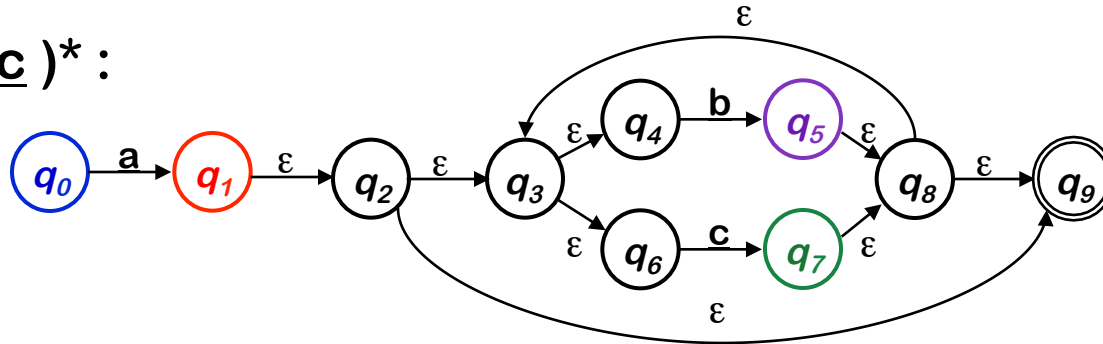
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
$s_2$	$q_5, q_8, q_9, q_3, q_4, q_6$			
$s_3$	$q_7, q_8, q_9, q_3, q_4, q_6$			

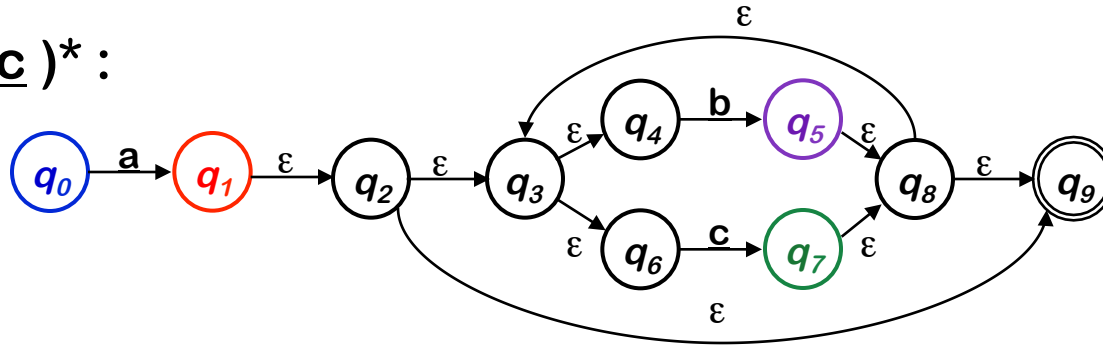
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
$s_2$	$q_5, q_8, q_9, q_3, q_4, q_6$	none		
$s_3$	$q_7, q_8, q_9, q_3, q_4, q_6$	none		

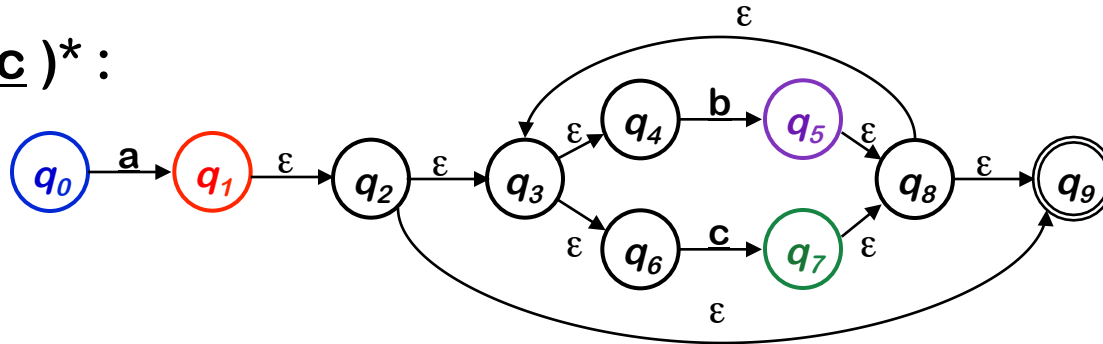
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
$s_2$	$q_5, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	
$s_3$	$q_7, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	

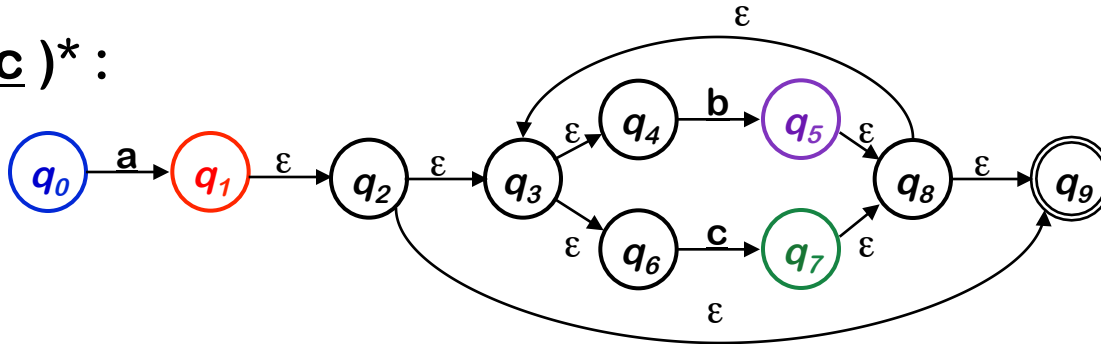
$a(b|c)^*$ :



Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
$s_2$	$q_5, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	$s_3$
$s_3$	$q_7, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	$s_3$

$a(b|c)^*$ :

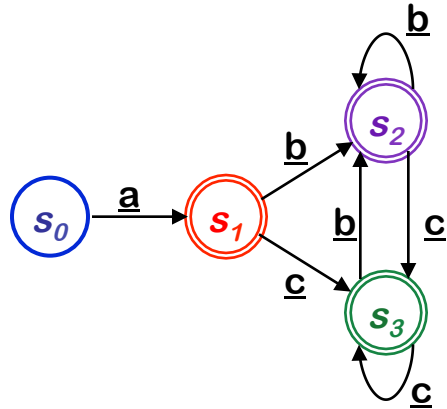


Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
$s_1$	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
$s_2$	$q_5, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	$s_3$
$s_3$	$q_7, q_8, q_9, q_3, q_4, q_6$	none	$s_2$	$s_3$

Final states

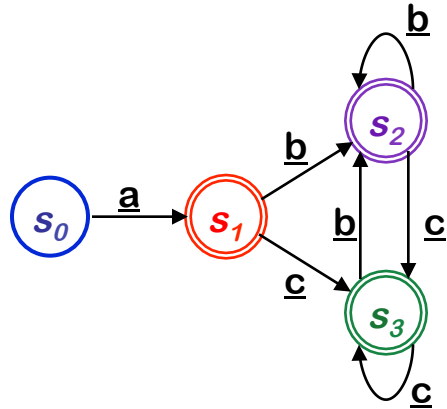
The DFA for  $\underline{a}(\underline{b} \mid \underline{c})^*$



$\delta$	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$s_1$	-	-
$s_1$	-	$s_2$	$s_3$
$s_2$	-	$s_2$	$s_3$
$s_3$	-	$s_2$	$s_3$

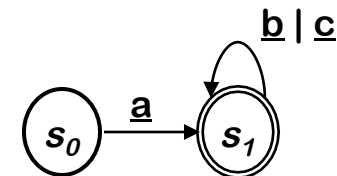
- Ends up smaller than the NFA
- All transitions are deterministic

The DFA for  $\underline{a}(\underline{b} \mid \underline{c})^*$



$\delta$	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$s_1$	-	-
$s_1$	-	$s_2$	$s_3$
$s_2$	-	$s_2$	$s_3$
$s_3$	-	$s_2$	$s_3$

- Ends up smaller than the NFA
- All transitions are deterministic
- Still not the smallest DFA



RE  $\rightarrow$  NFA (*Thompson's construction*)

- Build an NFA for each term
- Combine them with  $\epsilon$ -moves

NFA  $\rightarrow$  DFA (*subset construction*)

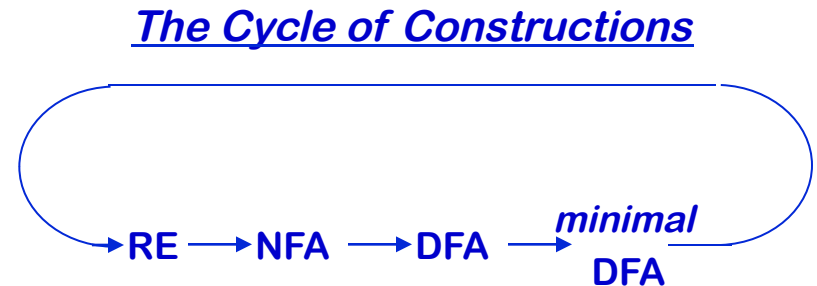
- Build the simulation

DFA  $\rightarrow$  Minimal DFA

- Hopcroft's algorithm

DFA  $\rightarrow$  RE (*not really part of scanner construction*)

- All pairs, all paths problem
- Union together paths from  $s_0$  to a final state



- Instruction Scheduling Project
  - Will be posted this Monday morning
  - Not a group project

# Syntax Analysis

Read EaC: 3.1 - 3.3