

# *CS415 Compilers*

## *Lexical Analysis*

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

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

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

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

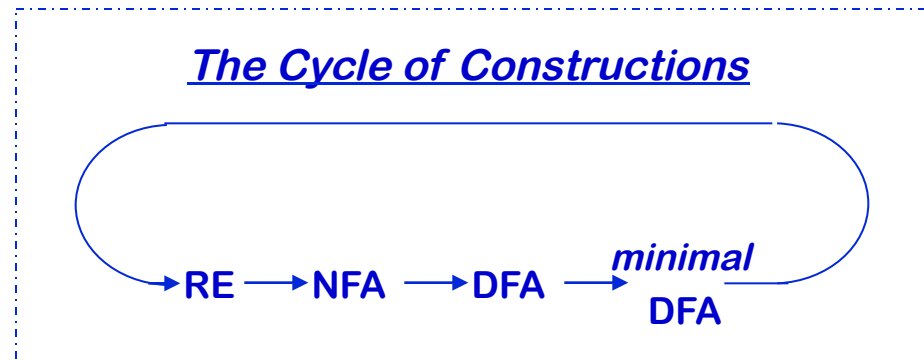
- Build the simulation

**DFA  $\rightarrow$  Minimal DFA**

- Hopcroft's algorithm

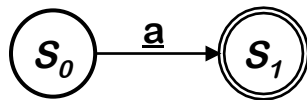
**DFA  $\rightarrow$  RE** (*Not part of the scanner construction*)

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



## Key idea

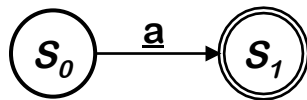
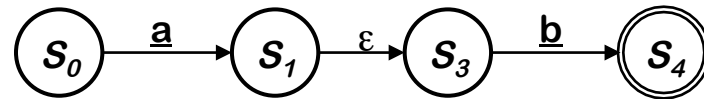
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for ab

Ken Thompson, CACM, 1968

## Key idea

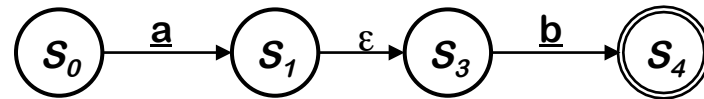
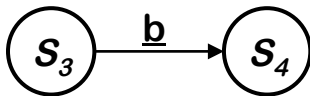
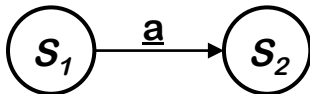
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for ab

Ken Thompson, CACM, 1968

## Key idea

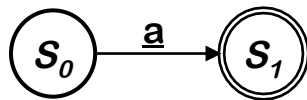
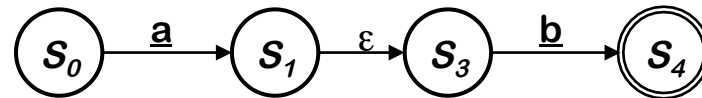
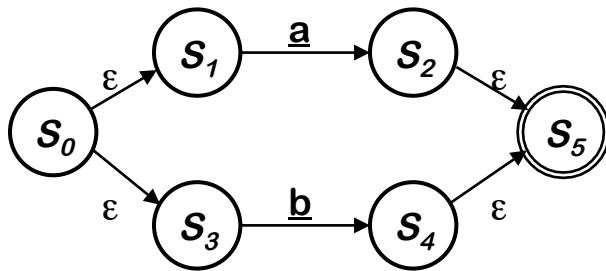
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

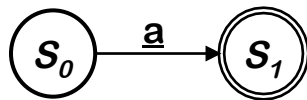
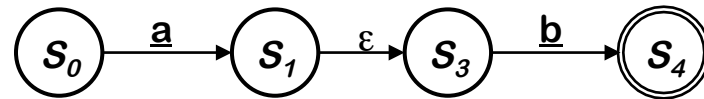
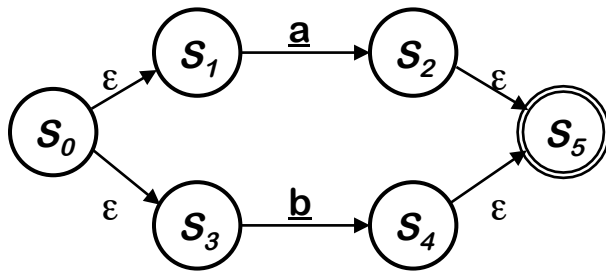
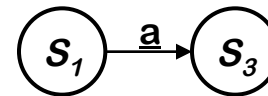
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

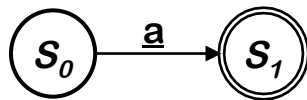
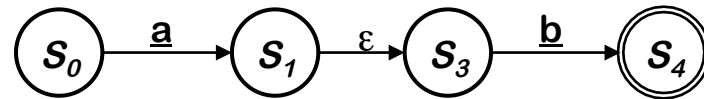
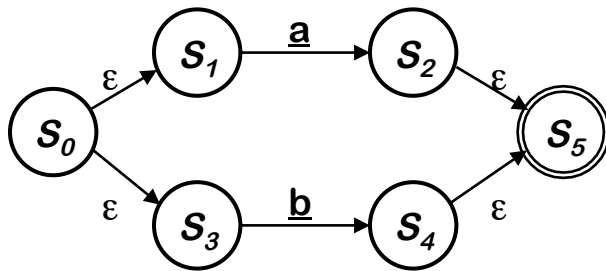
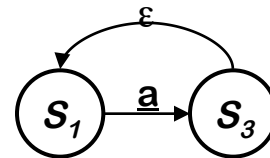
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for  $\underline{a}$ NFA for  $\underline{ab}$ NFA for  $\underline{a} \mid \underline{b}$ 

Ken Thompson, CACM, 1968

## Key idea

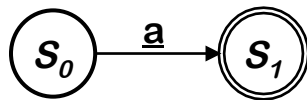
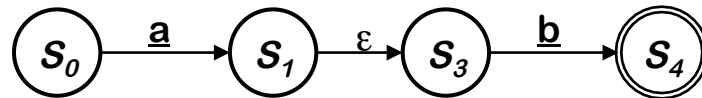
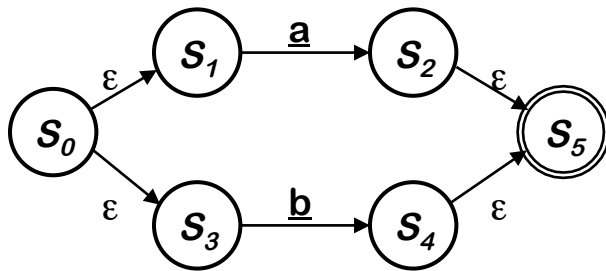
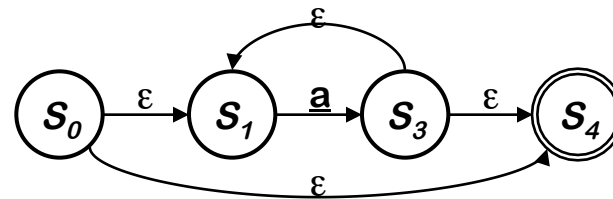
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | b

Ken Thompson, CACM, 1968

## Key idea

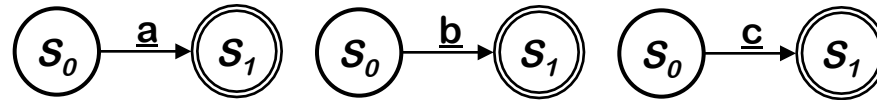
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
- Join them with  $\epsilon$  moves in precedence order

NFA for aNFA for abNFA for a | bNFA for a\*

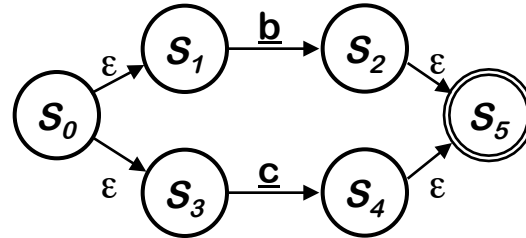
Ken Thompson, CACM, 1968

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

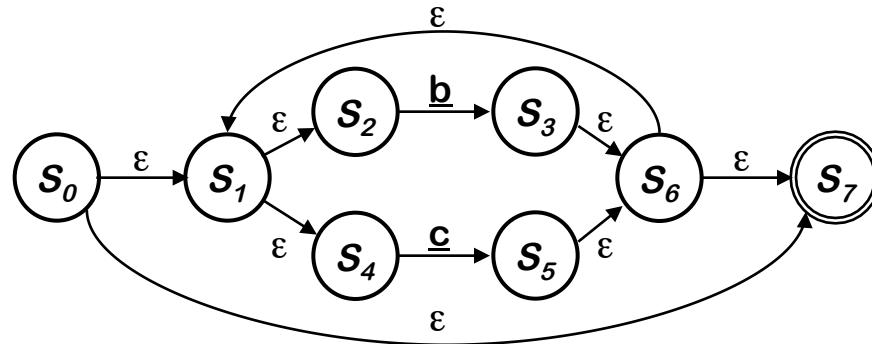
1.  $a, b, \& c$

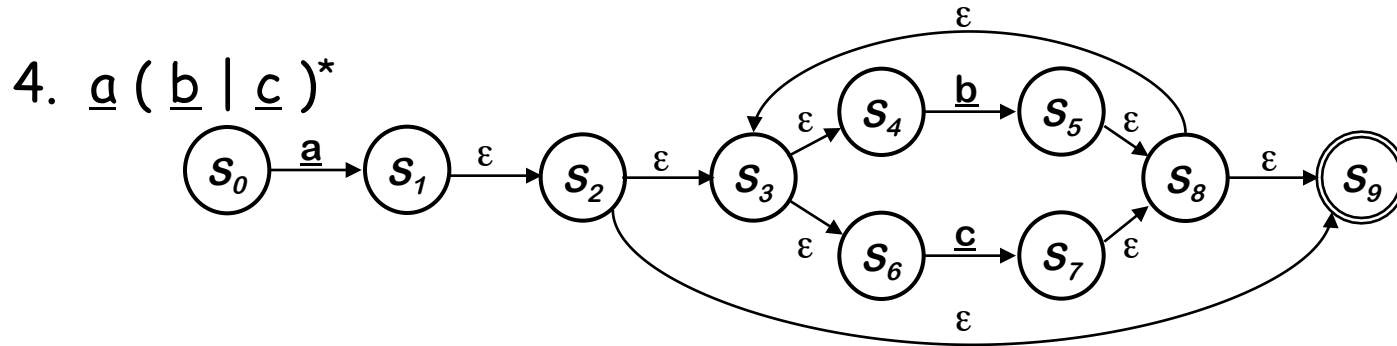


2.  $b|c$

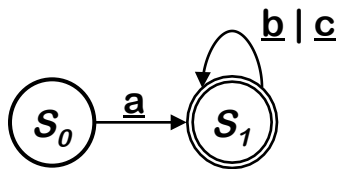


3.  $(b|c)^*$





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 (sketch):

- Start state derived from  $s_0$  of the NFA
- Take its  $\varepsilon$ -closure  $S_0 = \varepsilon$ -closure( $s_0$ )
- For each state  $S$ , compute  $move(S, \underline{a})$  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
         $T[s_i, a] \leftarrow s_i$ 

```

*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 symbol in each  $s_i$ .*

*It builds every possible NFA configuration.*

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

Example of a *fixed-point* computation

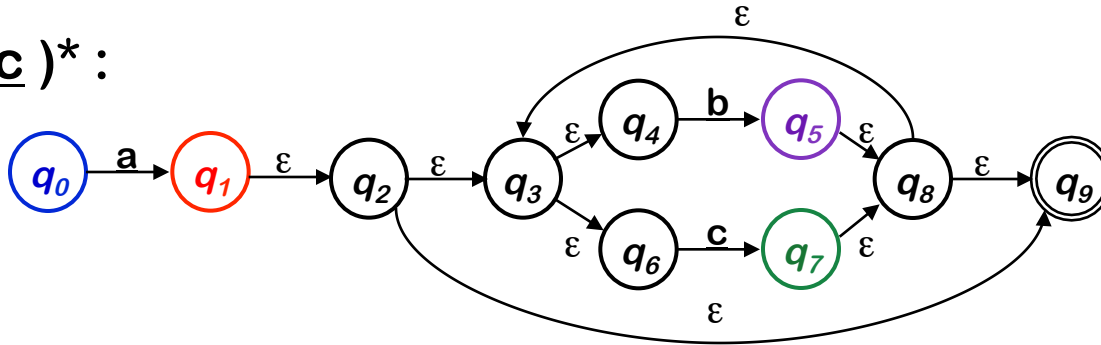
- Monotone construction of some finite set
- Halts when it stops adding to the set
- Proofs of halting & correctness are similar
- These computations arise in many contexts

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*

$a(b|c)^*$ :

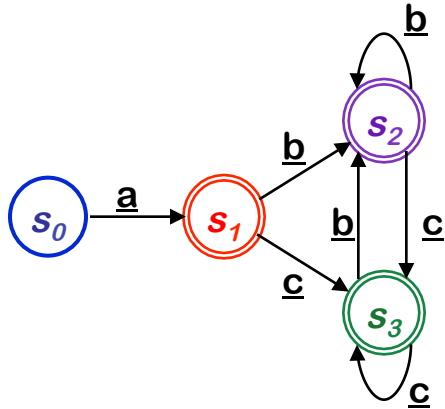


Applying the subset construction:

		$\epsilon$ -closure(move(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
$s_0$	$q_0$	$q_1, q_2, q_3, q_4, q_6, q_9$	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

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

The Cycle of Constructions



## The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

## The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

Two states are equivalent if and only if:

- $\forall a \in \Sigma$ , transitions on  $a$  lead to equivalent states (set) (DFA)

If  $a$ -transitions to distinct sets  $\Rightarrow$  states must be in distinct sets

## The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

Two states are equivalent if and only if:

- $\forall a \in \Sigma$ , transitions on  $a$  lead to equivalent states (DFA)

If  $a$ -transitions to distinct sets  $\Rightarrow$  states must be in distinct sets

A partition  $P$  of  $S$

- Each state  $s \in S$  is in exactly one set  $p_i \in P$
- The algorithm iteratively partitions the DFA's states

### Details of the algorithm

- Group states into maximal size sets, *optimistically*
- Iteratively subdivide those sets, as needed
- States that remain grouped together are equivalent

Initial partition,  $P_0$ , has two sets:  $\{F\}$  &  $\{Q-F\}$  ( $D = (Q, \Sigma, \delta, q_0, F)$ )

### Splitting a set (“partitioning a set by $\underline{a}$ ”)

- Assume  $q_a, \& q_b \in s$ , and  $\delta(q_a, \underline{a}) = q_x, \& \delta(q_b, \underline{a}) = q_y$
- If  $q_x \& q_y$  are not in the same set, then  $s$  must be split  
→  $q_a$  has transition on  $a$ ,  $q_b$  does not  $\Rightarrow \underline{a}$  splits  $s$

The algorithm

```

P ← { F, {Q-F} }
while ( P is still changing )
  T ← { }
  for each set S ∈ P
    T ← T ∪ split(S)
  P ← T

split(S):
  for each a ∈ Σ
    if a splits S into
      S1, S2, ... then
      return {S1, S2, ...}
  else return S

```

*This is a fixed-point algorithm!*

Why does this work?

- F is the set of final states
- Partition  $P \in 2^{|Q|}$
- Start off with 2 subsets of Q {F} and {Q-F}
- While loop takes  $P_i \rightarrow P_{i+1}$  by splitting 1 or more sets
- $P_{i+1}$  is at least one step closer to the partition with |Q| sets
- Maximum of |Q| splits

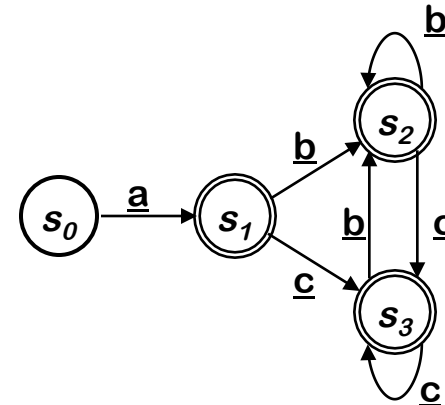
Note that

- Partitions are never combined

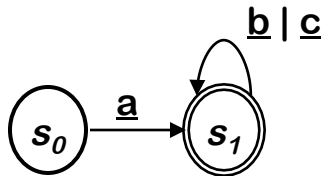
Then, apply the minimization algorithm

	Current Partition	Split on		
		<u>a</u>	<u>b</u>	<u>c</u>
$P_0$	$\{s_1, s_2, s_3\} \{s_0\}$	none	none	none

*final states*



To produce the minimal DFA



We observed that a human would design a simpler automaton than Thompson's construction & the subset construction did.

Minimizing that DFA produces the one that a human would design!

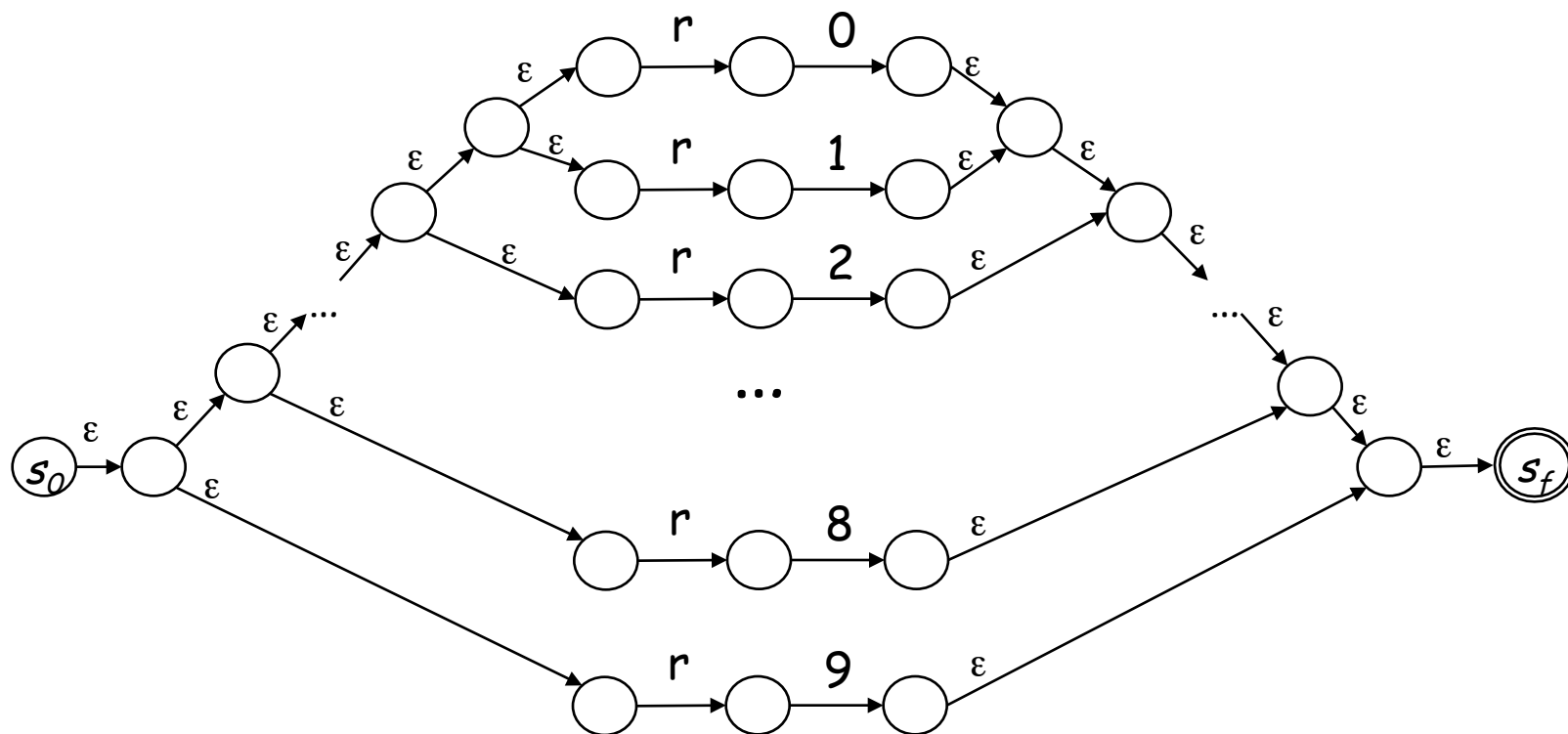
Start with a regular expression

r0 | r1 | r2 | r3 | r4 | r5 | r6 | r7 | r8 | r9

*The Cycle of Constructions*



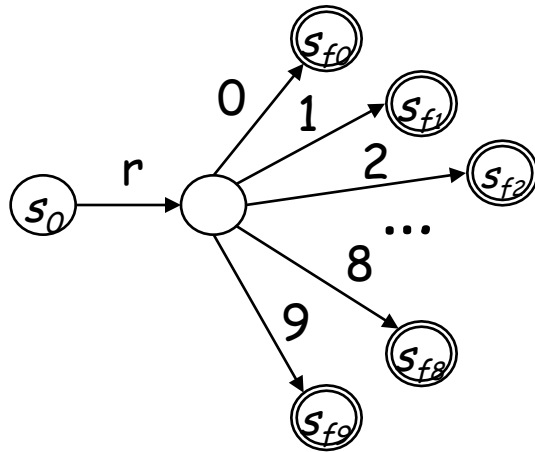
Thompson's construction produces



The Cycle of Constructions



The subset construction builds

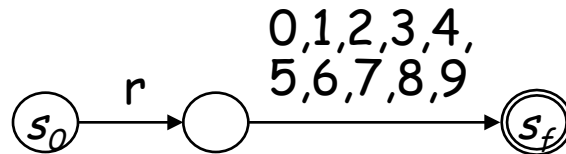


This is a DFA, but it has a lot of states ...

### The Cycle of Constructions



The DFA minimization algorithm builds



This looks like what a skilled compiler writer would do!

### The Cycle of Constructions



## Advantages of Regular Expressions

- Simple & powerful notation for specifying patterns
- Automatic construction of fast recognizers
- Many kinds of syntax can be specified with REs

Example — an expression grammar

$$\textit{Term} \rightarrow [a-zA-Z] ([a-zA-z] | [\underline{0}-\underline{9}])^*$$
$$\textit{Op} \rightarrow + | - | * | /$$
$$\textit{Expr} \rightarrow ( \textit{Term} \textit{Op} )^* \textit{Term}$$

Of course, this would generate a DFA ...

If REs are so useful ...

*Why not use them for everything?*

Not all languages are regular

$$RL's \subset CFL's \subset CSL's$$

You cannot construct DFA's to recognize these languages

- $L = \{p^k q^k\}$  *(parenthesis languages)*
- $L = \{wcw^r \mid w \in \Sigma^*\}$

Neither of these is a regular language *(nor an RE)*

But, this is a little subtle. You can construct DFA's for

- Strings with alternating 0's and 1's  
 $(\epsilon \mid 1)(01)^*(\epsilon \mid 0)$
- Strings with an even number of 0's and 1's
- Strings of bit patterns that represent binary numbers which are divisible by 5

1. Project 1 posted
2. Homework set 3 will be posted later today

**Next Class: Syntax Analysis**

Read EaC: 3.1 - 3.3