

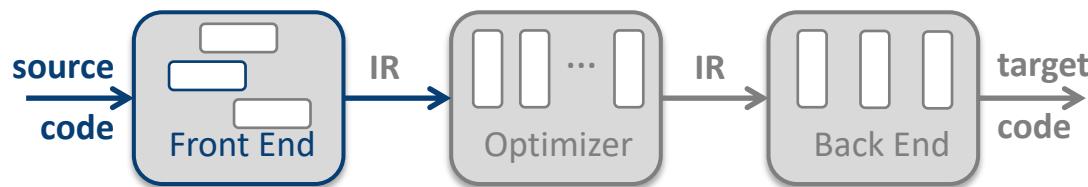


RICE

COMP 412  
FALL 2018

## Lexical Analysis, II

Comp 412



Copyright 2018, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.

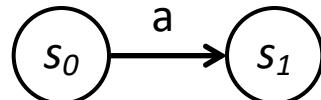
Chapter 2 in Eac2e



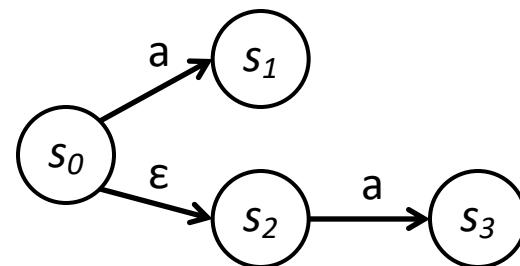
# Determinism (or not)

So far, we have only looked at deterministic automata, or DFAs

- DFA  $\cong$  Deterministic Finite Automaton
- Deterministic means that it has only one transition out of a state on a given character



*rather than*

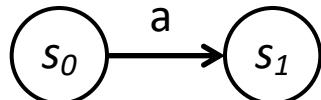




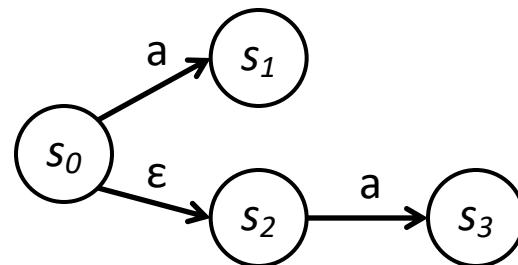
# Determinism (or not)

So far, we have only looked at deterministic automata, or DFAs

- DFA  $\approx$  Deterministic Finite Automaton
- Deterministic means that it has only one transition out of a state on a given character



*rather than*



- Can a finite automaton have multiple transitions out of a single state on the same character?
  - Yes, we call such an FA a Nondeterministic Finite Automaton
  - And, yes, the **NFA** is truly an odd notion ... but a useful one
- NFAs and DFAs are equivalent
  - The set of DFAs is a subset of the set of NFAs
  - For any NFA, we can build a DFA that simulates its behavior

⇒ We should not worry that the  $\epsilon$ -transitions do not consume input

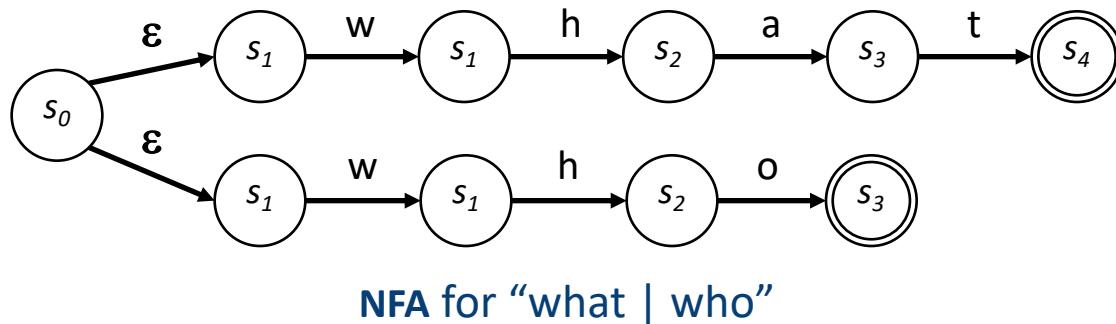


# Whoa. What Does That NFA “Mean”?

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  $x$ , ignoring  $\epsilon$ 's

Two models for NFA execution

1. To “run” the NFA, start in  $s_0$  and *guess* the right transition at each step <sup>†</sup>
2. To “run” the NFA, start in  $s_0$  and, at each non-deterministic choice, clone the NFA to pursue all possible paths. If any of the clones succeeds, *accept*



Note that this same operational definition works on a DFA

<sup>†</sup> See page 44 in EaC2e.

# Why Do We Care?

---



We need a construction that takes an RE to a DFA to a scanner.  
NFAs will help us get there.

## Overview:

1. Simple and direct construction of a **nondeterministic finite automaton (NFA)** to recognize a given RE
  - Easy to build in an algorithmic way
  - Key idea is to combine NFAs for the terms with  $\epsilon$ -transitions
2. Construct a **deterministic finite automaton (DFA)** that simulates the **NFA**
  - Use a set-of-states construction
3. Minimize the number of states in the **DFA**
  - We will look at 2 different algorithms: Hopcroft's & Brzozowski's
4. Generate the scanner code
  - Additional specifications needed for the actions

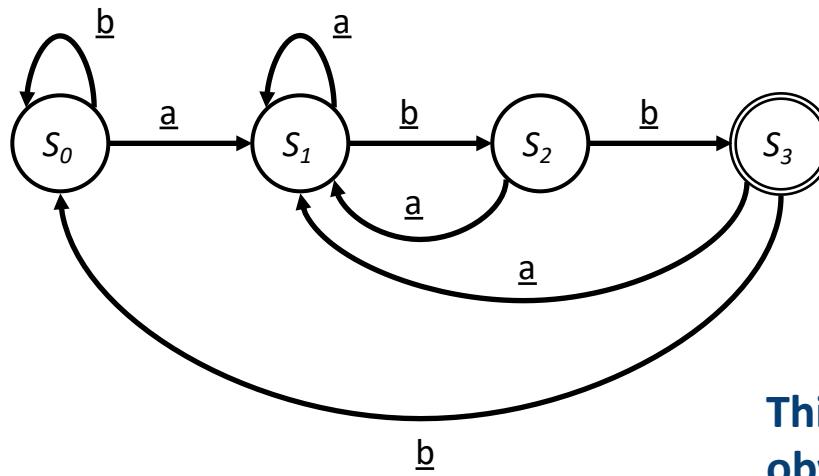
Optional, but worthwhile;  
reduces DFA size

**lex** and **flex** work along these lines



## Example of a DFA

Here is a DFA for  $(\underline{a} \mid \underline{b})^* \underline{abb}$



This DFA is not particularly obvious from the RE.

Each RE corresponds to one or more **deterministic finite automata** (DFAs)

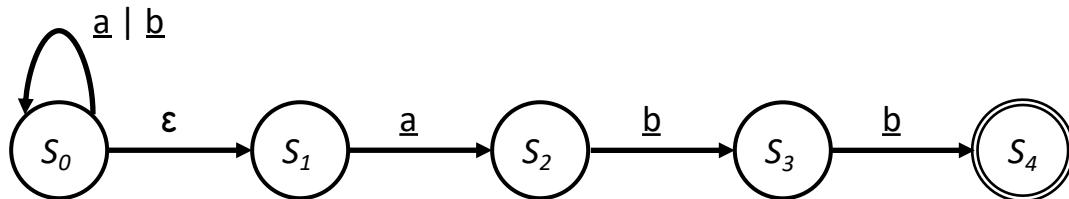
- We know a DFA exists for each RE
- The DFA may be hard to build directly
- Automatic techniques will build it for us ...

For algorithm aficionados in the class, this DFA is reminiscent of the way that the failure function works in the Knuth, Morris, & Pratt sub-linear time pattern matcher.



## Example as an NFA

Here is a simpler, more obvious NFA for  $(\underline{a} \mid \underline{b})^* \underline{abb}$



$(\underline{a} \mid \underline{b})^* \underline{abb}$

Here is an NFA for the same language

- The relationship between the **RE** and the **NFA** is more obvious
- The  $\epsilon$ -transition pastes together two **DFA**s to form a single **NFA**
- We can rewrite this **NFA** to eliminate the  $\epsilon$ -transition
  - $\epsilon$ -transitions are an odd and convenient quirk of **NFAs**
  - Eliminating this one makes it obvious that it has 2 transitions on  $\underline{a}$  from  $s_0$



# Relationship between NFAs and DFAs

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, but still true)*

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

⇒ NFA & DFA are equivalent in ability to recognize languages



# The Plan for Scanner Construction

**RE → NFA** (*Thompson's construction*)

- Build an **NFA** for each term in the **RE**
- Combine them in patterns that model the operators

**NFA → DFA** (*Subset construction*)

- Build a **DFA** that simulates the **NFA**

**DFA → Minimal DFA**

- Hopcroft's algorithm
- Brzozowski's algorithm

**Minimal DFA → Scanner**

- See § 2.5 in EaC2e

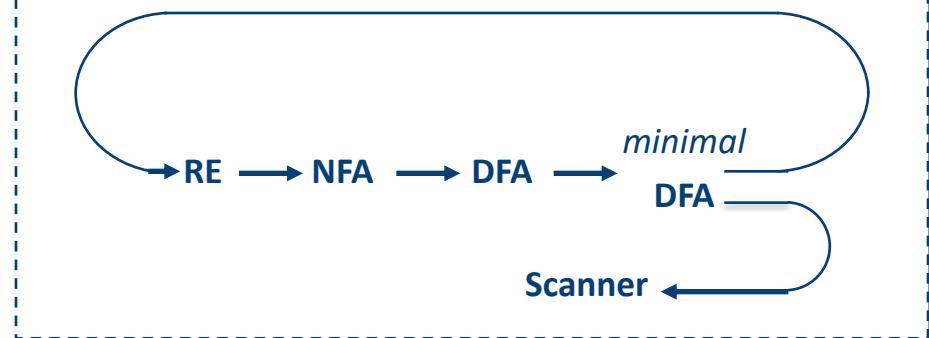
**DFA → RE**

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

## Automata Theory Moment

Taken together, the constructions in the cycle show that **REs**, **NFAs**, and **DFAs** are all equivalent in their expressive power.

### The Cycle of Constructions



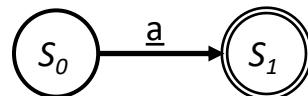
Taken together, these constructions prove that **DFAs** and **REs** are equivalent.

# RE $\rightarrow$ NFA using Thompson's Construction

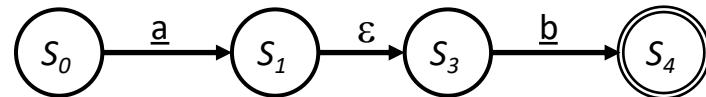


## Key idea

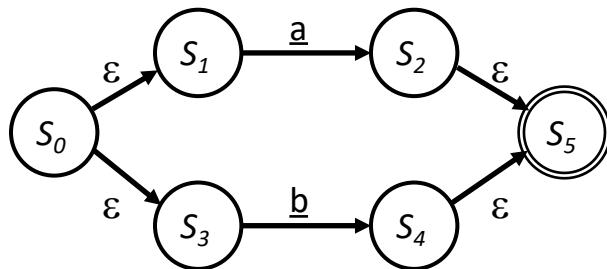
- For each symbol & each operator, we have an **NFA** pattern
- Join them with  $\epsilon$  moves in precedence order



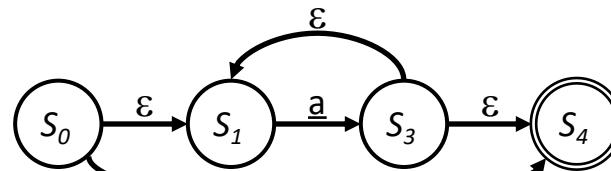
NFA for a



NFA for ab



NFA for a | b



NFA for a<sup>\*</sup>

### Precedence in REs:

Closure  
Concatenation  
Alternation

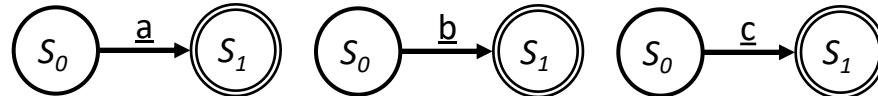
Ken Thompson, CACM, 1968



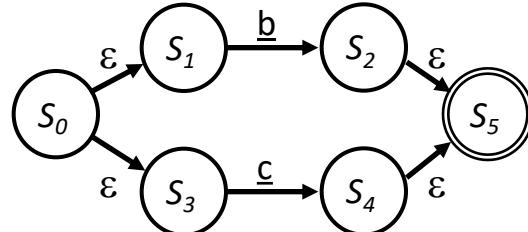
# Example of Thompson's Construction

Let's build an NFA for  $\underline{a} (\underline{b} \mid \underline{c})^*$

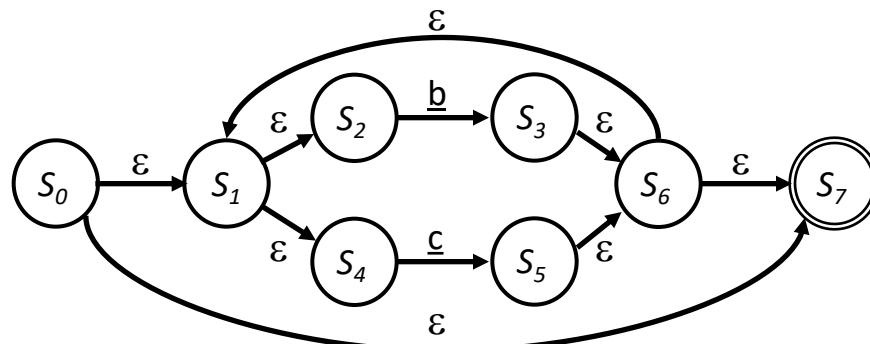
1.  $\underline{a}$ ,  $\underline{b}$ , &  $\underline{c}$



2.  $\underline{b} \mid \underline{c}$



3.  $(\underline{b} \mid \underline{c})^*$

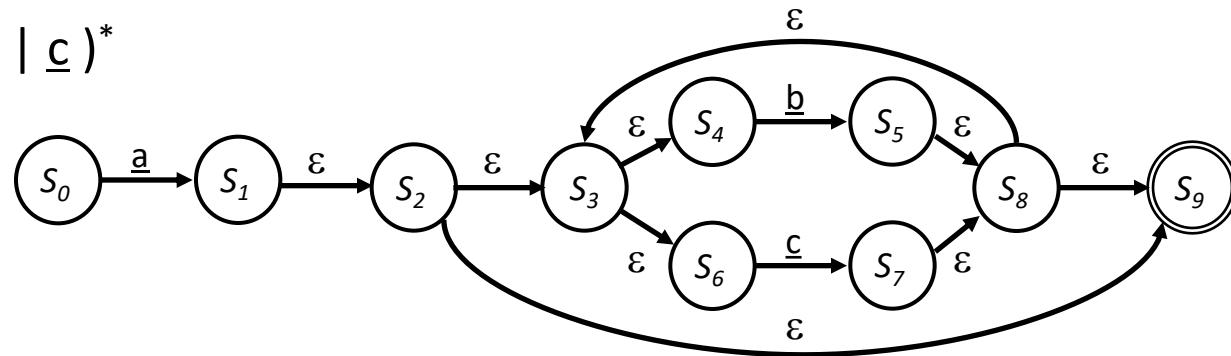


Note that states are being renamed at each step.

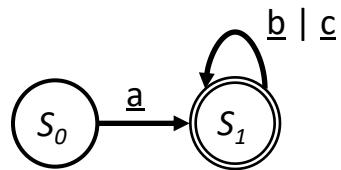
# Example of Thompson's Construction



4.  $\underline{a} (\underline{b} \mid \underline{c})^*$



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



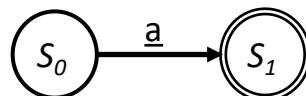
But, we can automate production of the more complex NFA version ...

# Thompson's Construction

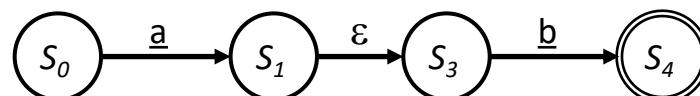


## Warning

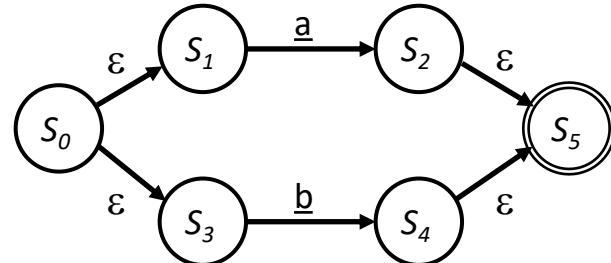
- You will be tempted to take shortcuts, such as leaving out some of the  $\epsilon$  transitions
- Do not do it. Memorize these four patterns. They will keep you out of trouble.



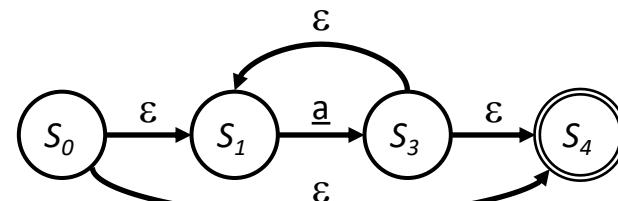
NFA for a



NFA for ab



NFA for a | b



NFA for a<sup>\*</sup>



# The Plan for Scanner Construction

**RE → NFA** (*Thompson's construction*)

- Build an **NFA** for each term in the **RE**
- Combine them in patterns that model the operators

**NFA → DFA** (*Subset construction*)

- Build a **DFA** that simulates the **NFA**

**DFA → Minimal DFA**

- Hopcroft's algorithm
- Brzozowski's algorithm

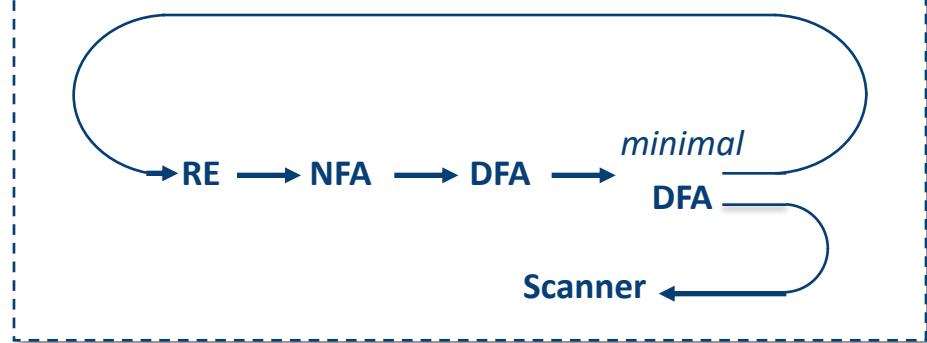
**Minimal DFA → Scanner**

- See § 2.5 in EaC2e

**DFA → RE**

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

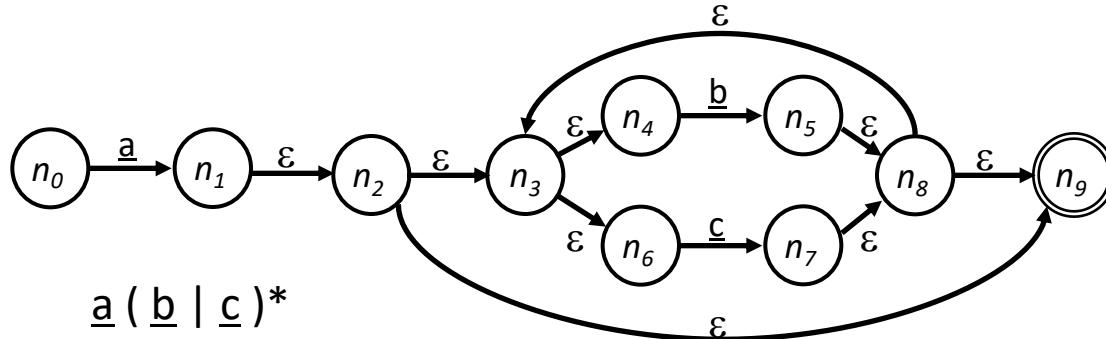
*The Cycle of Constructions*



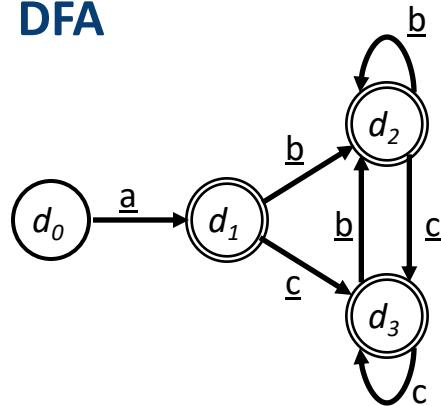


# Simulating an NFA with a DFA

NFA



DFA



Where the mapping between NFA states and DFA states is:

DFA	NFA
$d_0$	$n_0$
$d_1$	$n_1 \ n_2 \ n_3 \ n_4 \ n_6 \ n_9$
$d_2$	$n_5 \ n_8 \ n_9 \ n_3 \ n_4 \ n_6$
$d_3$	$n_7 \ n_8 \ n_9 \ n_3 \ n_4 \ n_6$



# NFA → DFA with Subset Construction

The subset construction builds a DFA that simulates the NFA

## Two key functions

- $\text{Move}(s_i, \underline{a})$  is the set of states reachable from  $s_i$  by  $\underline{a}$
- $\text{FollowEpsilon}(s_i)$  is the set of states reachable from  $s_i$  by  $\epsilon$

## The algorithm

- Derive the DFA's start state from  $n_0$  of the NFA
- Add all states reachable from  $n_0$  by following  $\epsilon$ 
  - $d_0 = \text{FollowEpsilon}(\{n_0\})$
  - Let  $D = \{ d_0 \}$
- For  $\alpha \in \Sigma$ , compute  $\text{FollowEpsilon}(\text{Move}(d_0, \alpha))$ 
  - If this creates a new state, add it to  $D$
- Iterate until no more states are added

It sounds more complex than it is...

Rabin & Scott, 1959

15

Any DFA state that contains a final state of the NFA becomes a final state of the DFA.



## NFA → DFA with Subset Construction

### The algorithm:

```
 $d_0 \leftarrow \text{FollowEpsilon}(\{n_0\})$ 
 $D \leftarrow \{d_0\}$ 
 $W \leftarrow \{d_0\}$ 
while ( $W \neq \emptyset$ ) {
    select and remove  $s$  from  $W$ 
    for each  $\alpha \in \Sigma$  {
         $t \leftarrow \text{FollowEpsilon}(\text{Move}(s, \alpha))$ 
         $T[s, \alpha] \leftarrow t$ 
        if ( $t \notin D$ ) then {
            add  $t$  to  $D$ 
            add  $t$  to  $W$ 
        }
    }
}
```

$d_0$  is a set of states  
D & W are sets of sets of states

### The algorithm halts:

1.  $D$  contains no duplicates (test before addition)
2.  $2^{\{\text{NFA states}\}}$  is finite
3. while loop adds to  $D$ , but does not remove from  $D$  (monotone)

⇒ the loop halts

$D$  contains all the reachable NFA states

*It tries each character in each  $d_i$ .*

*It builds every possible NFA configuration.*

⇒  $D$  and  $T$  form the DFA

This test is a little tricky

# NFA → DFA with Subset Construction

---



## 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

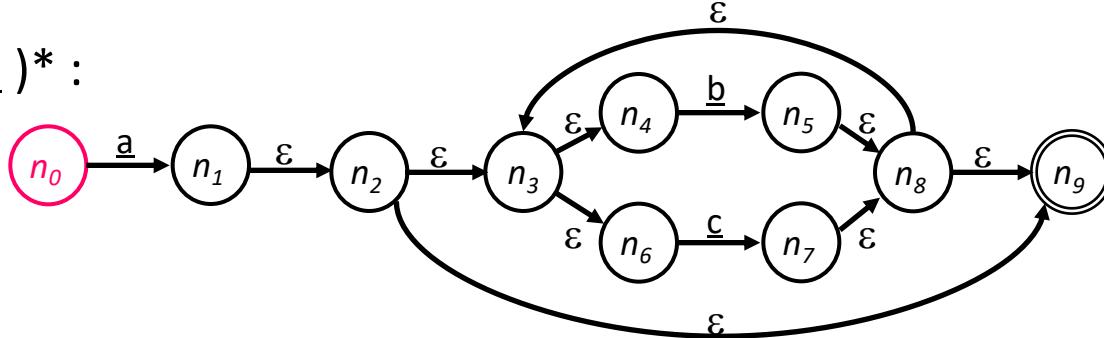
- Canonical construction of sets of LR(1) items
  - Quite similar to the subset construction
- Classic data-flow analysis & Gaussian Elimination
  - Solving sets of simultaneous set equations

*We will see many more fixed-point computations*



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

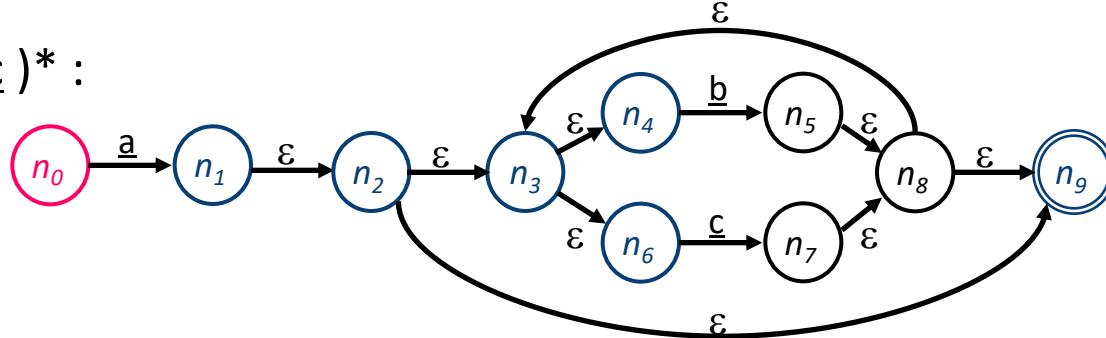


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	a	<u>b</u>	<u>c</u>
$d_0$	$n_0$			



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

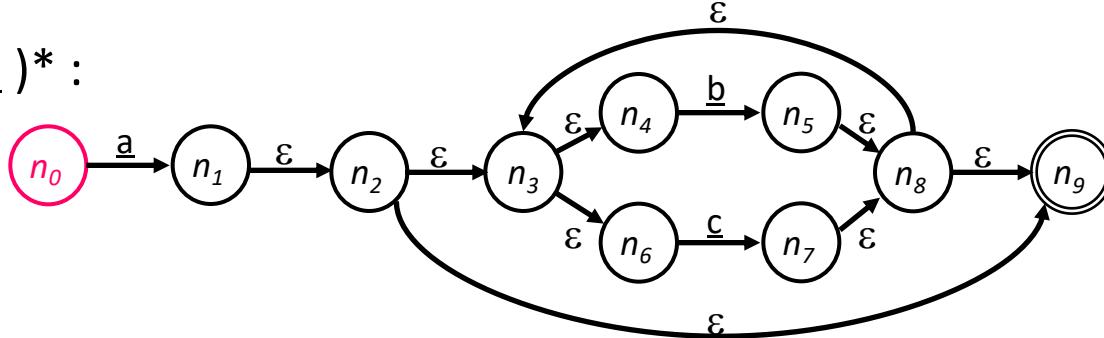


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$		



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

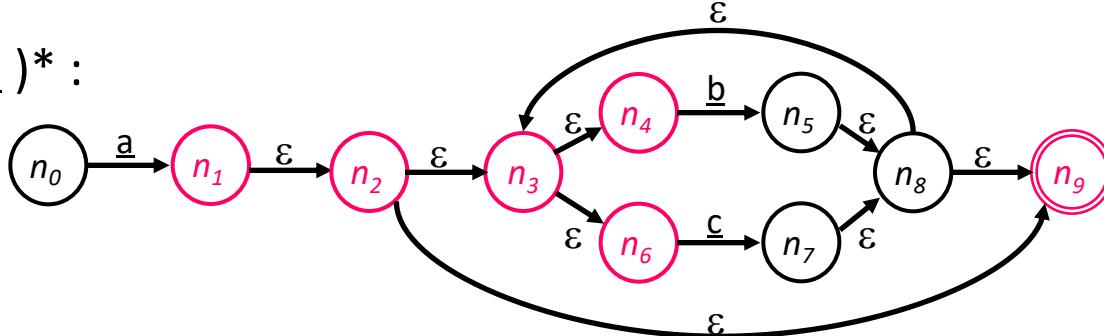


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	none



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

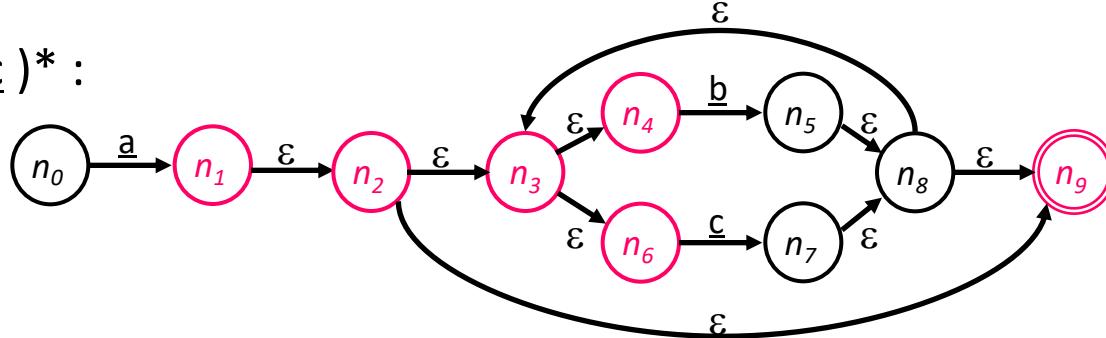


States		FollowEpsilon ( Move( s, * ) )		
DFA	NFA	a	b	c
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	none



# NFA → DFA with Subset Construction

a ( b | c )\* :

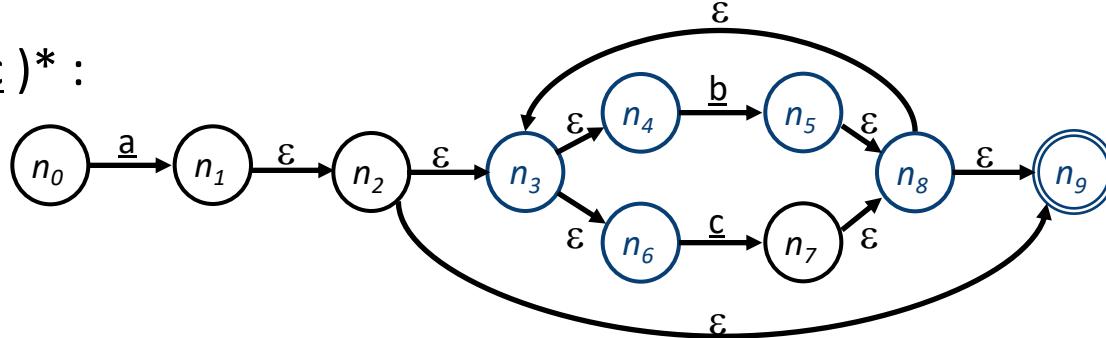


States		FollowEpsilon ( Move( s, * ) )		
DFA	NFA	a	b	c
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	none
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none		



# NFA → DFA with Subset Construction

a ( b | c )\* :

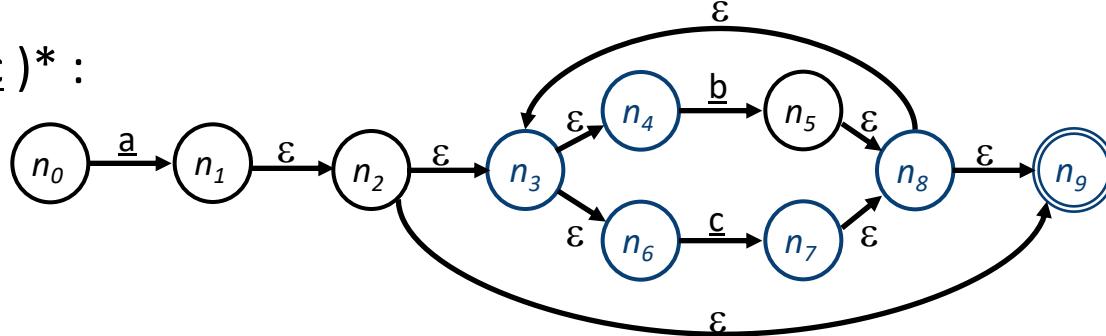


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	a	b	c
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

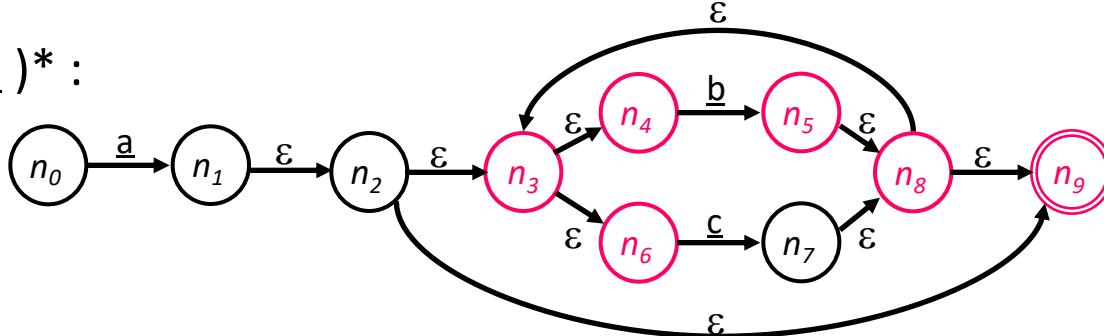


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

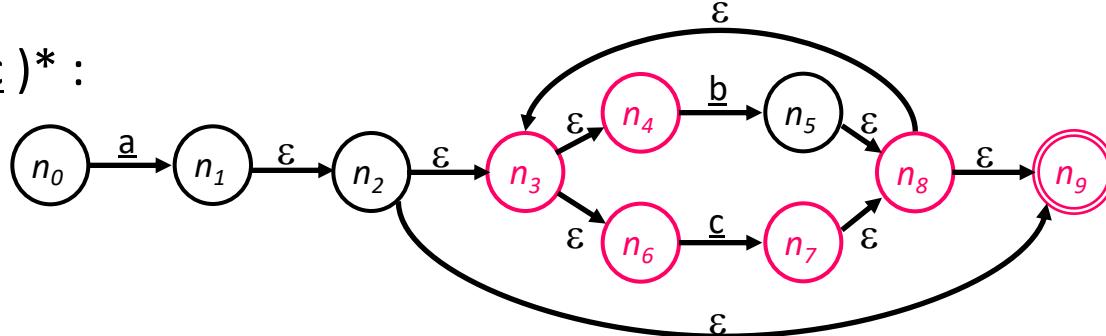


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$			



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

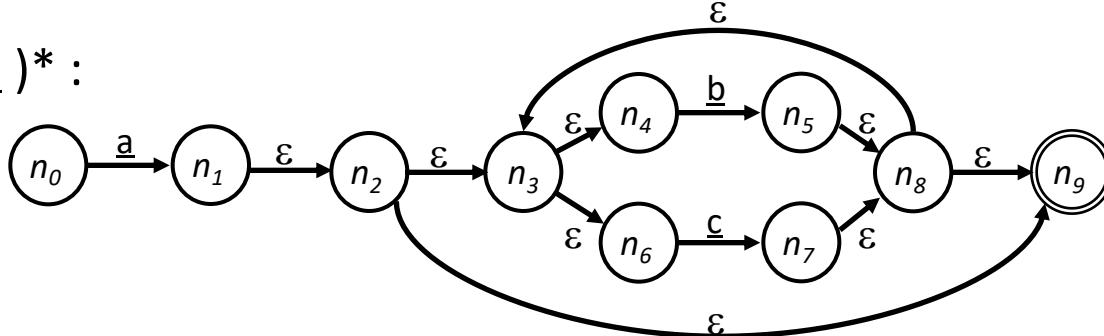


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$			
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$			



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

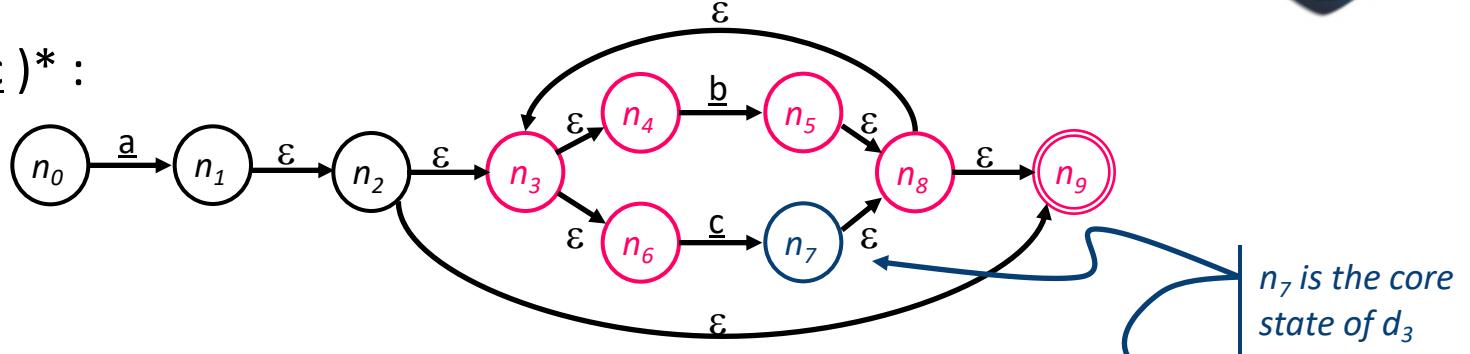


States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>		
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>		



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :

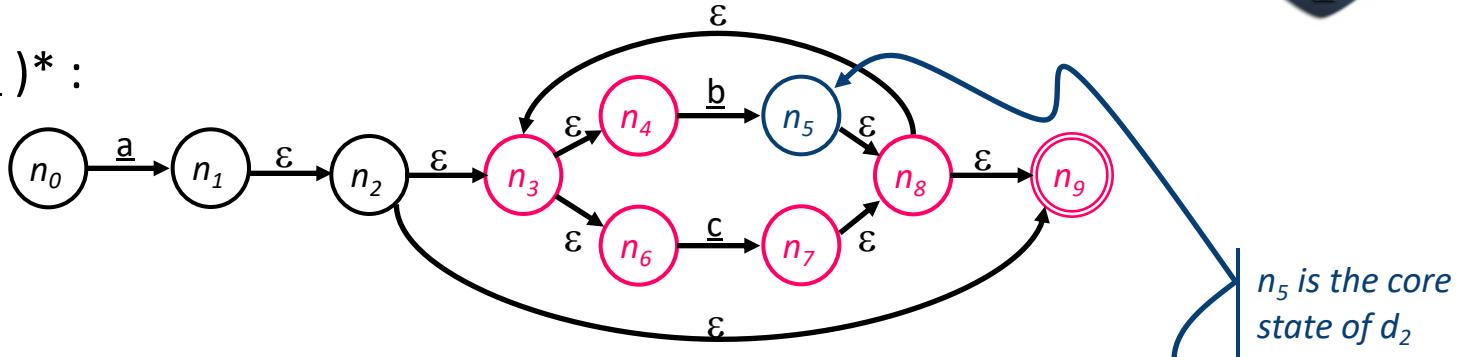


States		<i>FollowEpsilon (Move(s, *))</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	none
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	none	$d_2$	$d_3$
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	none		

# NFA → DFA with Subset Construction



a ( b | c )\* :

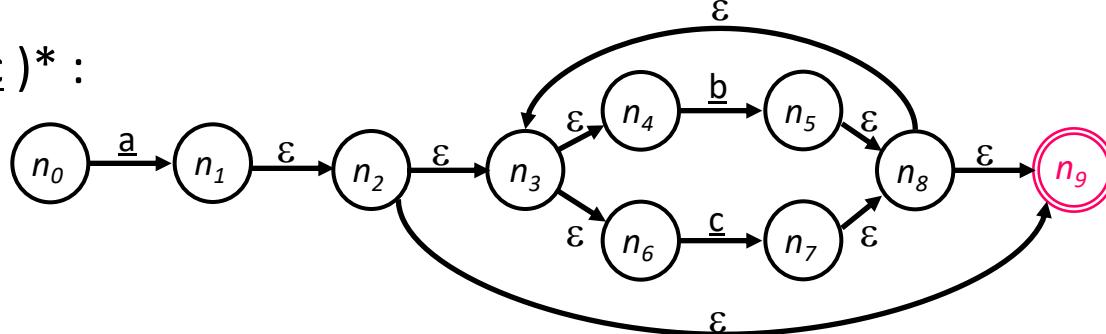


States		FollowEpsilon ( Move( s, * ) )		
DFA	NFA	a	b	c
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	none
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	none	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	none	$d_2$	$d_3$
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	none	$d_2$	$d_3$



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :



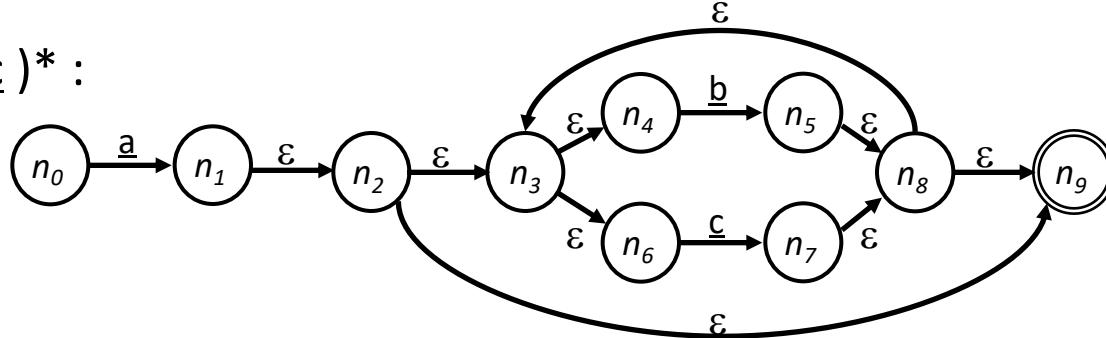
States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>	$d_2$	$d_3$
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>	$d_2$	$d_3$

Final states because of  $n_9$



# NFA → DFA with Subset Construction

$a (\underline{b} \mid \underline{c})^*$ :



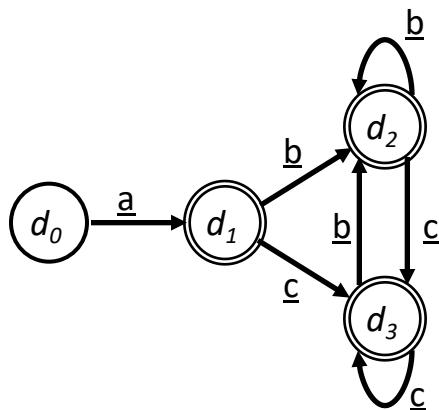
States		<i>FollowEpsilon ( Move( s, * ) )</i>		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
$d_0$	$n_0$	$d_1$	<i>none</i>	<i>none</i>
$d_1$	$n_1 \ n_2 \ n_3$ $n_4 \ n_6 \ n_9$	<i>none</i>	$d_2$	$d_3$
$d_2$	$n_5 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>	$d_2$	$d_3$
$d_3$	$n_7 \ n_8 \ n_9$ $n_3 \ n_4 \ n_6$	<i>none</i>	$d_2$	$d_3$

Transition table for the DFA



# NFA → DFA with Subset Construction

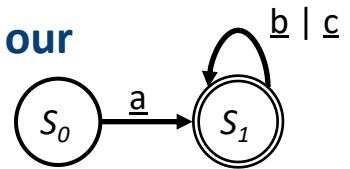
The DFA for  $\underline{a} (\underline{b} \mid \underline{c})^*$



	<u>a</u>	<u>b</u>	<u>c</u>
<u>d</u> <sub>0</sub>	<u>d</u> <sub>1</sub>	none	none
<u>d</u> <sub>1</sub>	none	<u>d</u> <sub>2</sub>	<u>d</u> <sub>3</sub>
<u>d</u> <sub>2</sub>	none	<u>d</u> <sub>2</sub>	<u>d</u> <sub>3</sub>
<u>d</u> <sub>3</sub>	none	<u>d</u> <sub>2</sub>	<u>d</u> <sub>3</sub>

- Much smaller than the NFA (no  $\epsilon$ -transitions)
- All transitions are deterministic
- Use same code skeleton as before

But, remember, our goal was:



Rabin and Scott, 1959 (page 8)

chines are more general than the ordinary ones, but this is not the case. We shall give a direct construction of an ordinary automaton, defining exactly the same set of tapes as a given nondeterministic machine.

**Definition 11.** Let  $\mathfrak{A} = (S, M, S_0, F)$  be a nondeterministic automaton.  $\mathfrak{D}(\mathfrak{A})$  is the system  $(T, N, t_0, G)$  where  $T$  is the set of all subsets of  $S$ ,  $N$  is a function on  $T \times \Sigma$  such that  $N(t, \sigma)$  is the union of the sets  $M(s, \sigma)$  for  $s \in t$ ,  $t_0 = S_0$ , and  $G$  is the set of all subsets of  $S$  containing at least one member of  $F$ .

→ Clearly  $\mathfrak{D}(\mathfrak{A})$  is an ordinary automaton, but it is actually equivalent to  $\mathfrak{A}$ .

**Theorem 11.** If  $\mathfrak{A}$  is a nondeterministic automaton, then  $T(\mathfrak{A}) = T(\mathcal{D}(\mathfrak{A}))$ .

*Proof:* Assume first that  $s_0, s_1, \dots$  states satisfying the conditions of Definition 10. we show by induction that for  $k \leq n$ ,  $s_k$  is in  $N(t_0, x_k)$ . For  $k=0$ ,  $N(t_0, x_k) = N(t_0, \Delta) = t_0 = S_0$  and we were given that  $s_0$  is in  $S_0$ . Assume the result for  $k-1$ . By definition,  $N(t_0, x_k) = N(N(t_0, x_{k-1}), \sigma_{k-1})$ . But we have assumed  $s_{k-1}$  is in  $N(t_0, x_{k-1})$  so that from the definition of  $N$  we have  $M(s_{k-1}, \sigma_{k-1}) \subset N(t_0, x_k)$ . However,  $s_k$  is in  $M(s_{k-1}, \sigma_{k-1})$ , and so the result is established. In particular  $s_n$  is in  $N(t_0, x_n) = N(t_0, x)$ , and since  $s_n$  is in  $F$ , we have  $N(t_0, x)$  in  $G$ , which proves that  $x$  is in  $T(\mathcal{D}(\mathfrak{A}))$ . Hence, we have shown that

$$T(\mathfrak{A}) \subset T(\mathfrak{D}(\mathfrak{A})).$$

Assume next that a tape  $x = \sigma_0\sigma_1\dots\sigma_{n-1}$  is in  $T(\mathfrak{D}(\mathfrak{A}))$ . Let for each  $k \leq n$ ,  $t_k = N(t_0, x_k)$ . We shall work backwards. First, we know that  $t_n$  is in  $G$ . Let then  $s_n$  be any internal state of  $\mathfrak{A}$  such that  $s_n$  is in  $t_n$  and  $s_n$  is in  $F$ . Since  $s_n$  is in

$$t_n = N(t_{n-1}, \sigma_{n-1}),$$

we have from the definition of  $N$  that  $s_n$  is in  $M(s_{n-1}; \sigma_{n-1})$  for some  $s_{n-1}$  in  $t_{n-1}$ . But

**Definition 12.** Let  $\mathfrak{A} = (S, M, S_0, F)$  be a nondeterministic automaton. The dual of  $\mathfrak{A}$  is the machine  $\mathfrak{A}^* = (S, M^*, F, S_0)$  where the function  $M^*$  is defined by the condition

$s'$  is in  $M^*(s, \sigma)$  if and only if  $s$  is in  $M(s', \sigma)$ .

Notice that we have at once the equation  $\mathfrak{A}^{**} = \mathfrak{A}$ . The relation between the sets defined by an automaton and its dual is as follows.

**Theorem 12.** If  $\mathfrak{A}$  is a nondeterministic automaton, then  $T(\mathfrak{A}^*) = T(\mathfrak{A})^*$ .

*Proof:* In view of the equality  $\mathfrak{A}^{**} = \mathfrak{A}$ , we need only show  $T(\mathfrak{A}^*) \subset T(\mathfrak{A})^*$ . Let  $x = \sigma_0\sigma_1 \dots \sigma_{n-1}$  be a tape of  $T(\mathfrak{A}^*)$ . We shall show that  $x^*$  is in  $T(\mathfrak{A})$ . Let  $s_0, s_1, \dots, s_n$  be a sequence of internal states of  $\mathfrak{A}^*$  such that  $s_0$  and  $s_k$  is in  $M^*(s_{k-1}, \sigma_{k-1})$  for  $k=1, 2, \dots, n$ . Define a new sequence  $s'_0, s'_1, \dots, s'_n$  by the equation  $s'_k = s_{n-k}$  for  $k \leq n$ . Obviously,  $s'_0$  is in  $S_0$  and  $s'_n$  is in  $F$ . Further, for  $k > 0$  and  $k \leq n$ ,  $s'_{k-1} = s_{n-k+1}$  is in  $M^*(s_{n-k}, \sigma_{n-k})$ , or in other words,  $s_{n-k} = s'_k$  is in  $M(s'_{k-1}, \sigma_{n-k})$ . Now defining a new sequence of symbols  $\sigma'_0\sigma'_1 \dots \sigma'_{n-1}$  by the formula  $\sigma'_k = \sigma_{n-k-1}$ , we see that  $\sigma'_{k-1} = \sigma_{n-k}$  and  $\sigma'_0\sigma'_1 \dots \sigma'_{n-1} = x^*$ . Thus,  $x^*$  is in  $T(\mathfrak{A})$  as was to be proved.

It should be noted that Theorem 12 together with Theorem 11 yields a direct construction and proof for Theorem 4 of Section 3 which was first proved by the indirect method of Theorem 1. In the next section we make heavy use of the direct constructions supplied by the nondeterministic machines to obtain results not easily apparent from the mathematical characterizations of Theorems 1 and 2.

## 6. Further closure properties

Simplifying a result due originally to Kleene, Myhill in unpublished work has shown that the class  $\mathcal{T}$  can be characterized as the least class of sets of tapes containing the finite sets and closed under some simple operations on sets of tapes. We indicate here a different proof using



# The Plan for Scanner Construction

**RE → NFA** (*Thompson's construction*)

- Build an **NFA** for each term in the **RE**
- Combine them in patterns that model the operators

**NFA → DFA** (*Subset construction*)

- Build a **DFA** that simulates the **NFA**

**DFA → Minimal DFA**

- Hopcroft's algorithm
- Brzozowski's algorithm

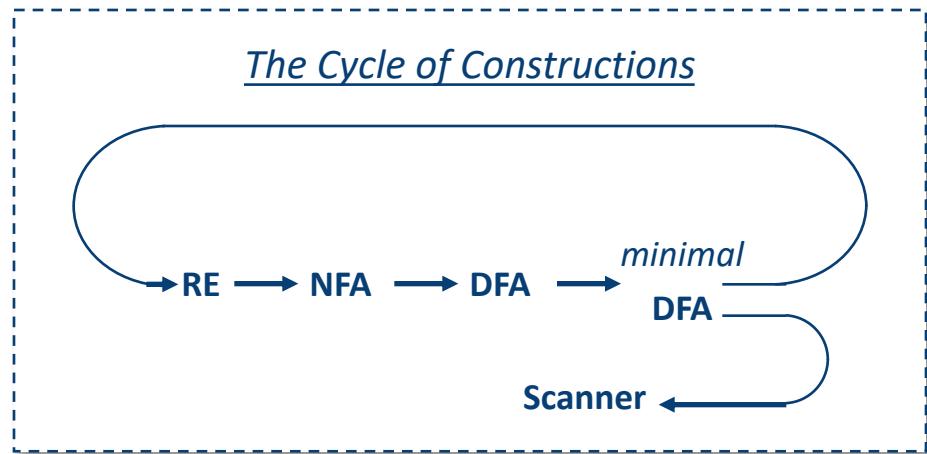
**Minimal DFA → Scanner**

- See § 2.5 in EaC2e

**DFA → RE**

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

*The Cycle of Constructions*

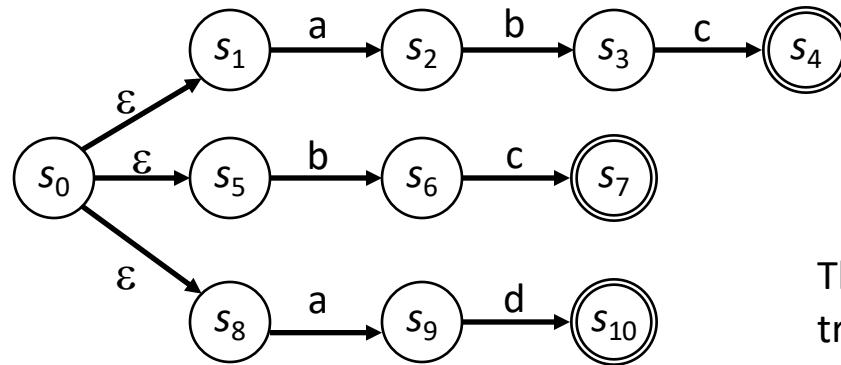


# Brzozowski's Algorithm for DFA Minimization



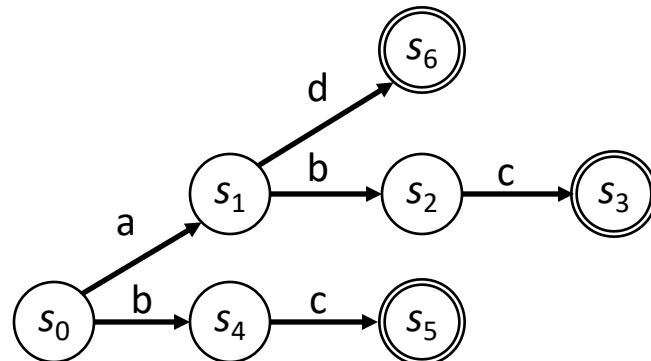
## The Intuition

- The subset construction merges prefixes in the NFA



abc | bc | ad

Thompson's construction would leave  $\epsilon$ -transitions between each single-character automaton



Subset construction eliminates  $\epsilon$ -transitions and merges the paths for a. It leaves duplicate tails, such as bc, intact.



# Brzozowski's Algorithm

## Idea: Use The Subset Construction Twice

- For an **NFA**  $N$ 
  - Let  $\text{reverse}(N)$  be the **NFA** constructed by making initial state final, adding a new start state with an  $\epsilon$ -transition to each previously final state, and reversing the other edges
  - Let  $\text{subset}(N)$  be the **DFA** produced by the subset construction on  $N$
  - Let  $\text{reachable}(N)$  be  $N$  after removing any states that are not reachable from the initial state
- Then,

$\text{reachable}(\text{subset}(\text{reverse}(\text{reachable}(\text{subset}(\text{reverse}(N)))))$ )

is a minimal **DFA** that implements  $N$  [Brzozowski, 1962]

***Not everyone finds this result to be intuitive.***

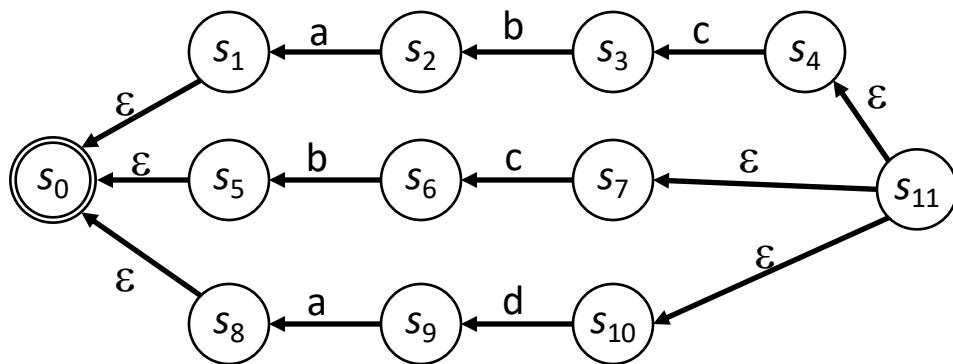
***Neither algorithm dominates the other.***



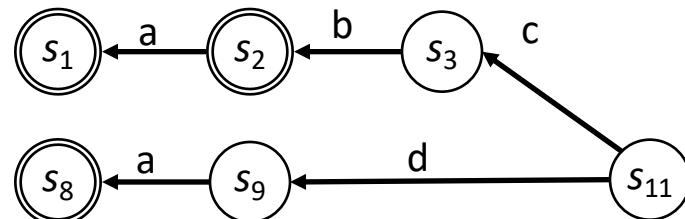
# Brzozowski's Algorithm

## Step 1

- The subset construction on  $\text{reverse}(\text{NFA})$  merges suffixes in original NFA



Reversed NFA



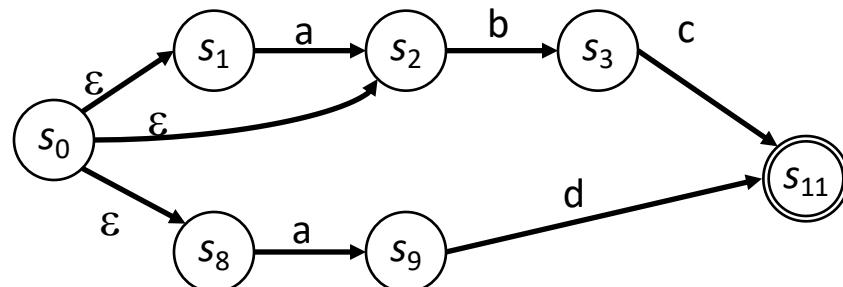
subset( $\text{reverse}(\text{NFA})$ )

# Brzozowski's Algorithm

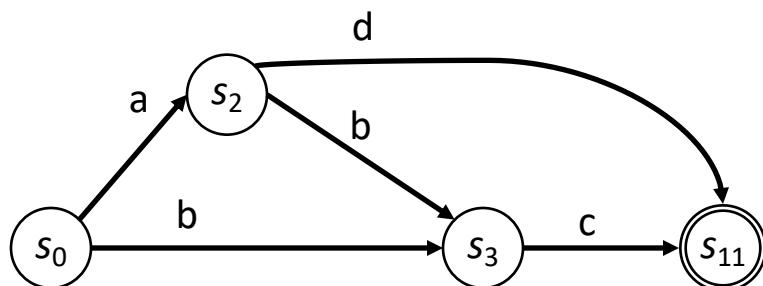


## Step 2

- Reverse it again & use subset to merge prefixes ...



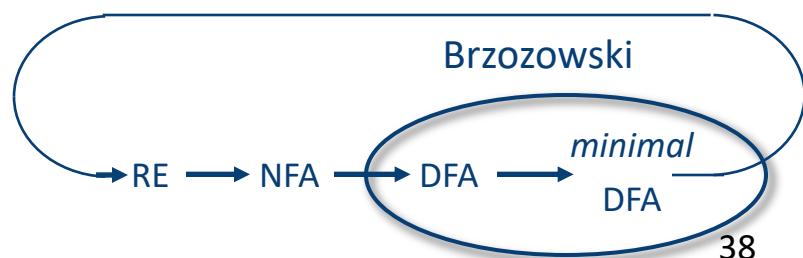
Reverse it, again



And subset it, again

Minimal DFA

The Cycle of Constructions



# Abbreviated Register Specification



**Start with a regular expression**

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

Register names from zero to nine

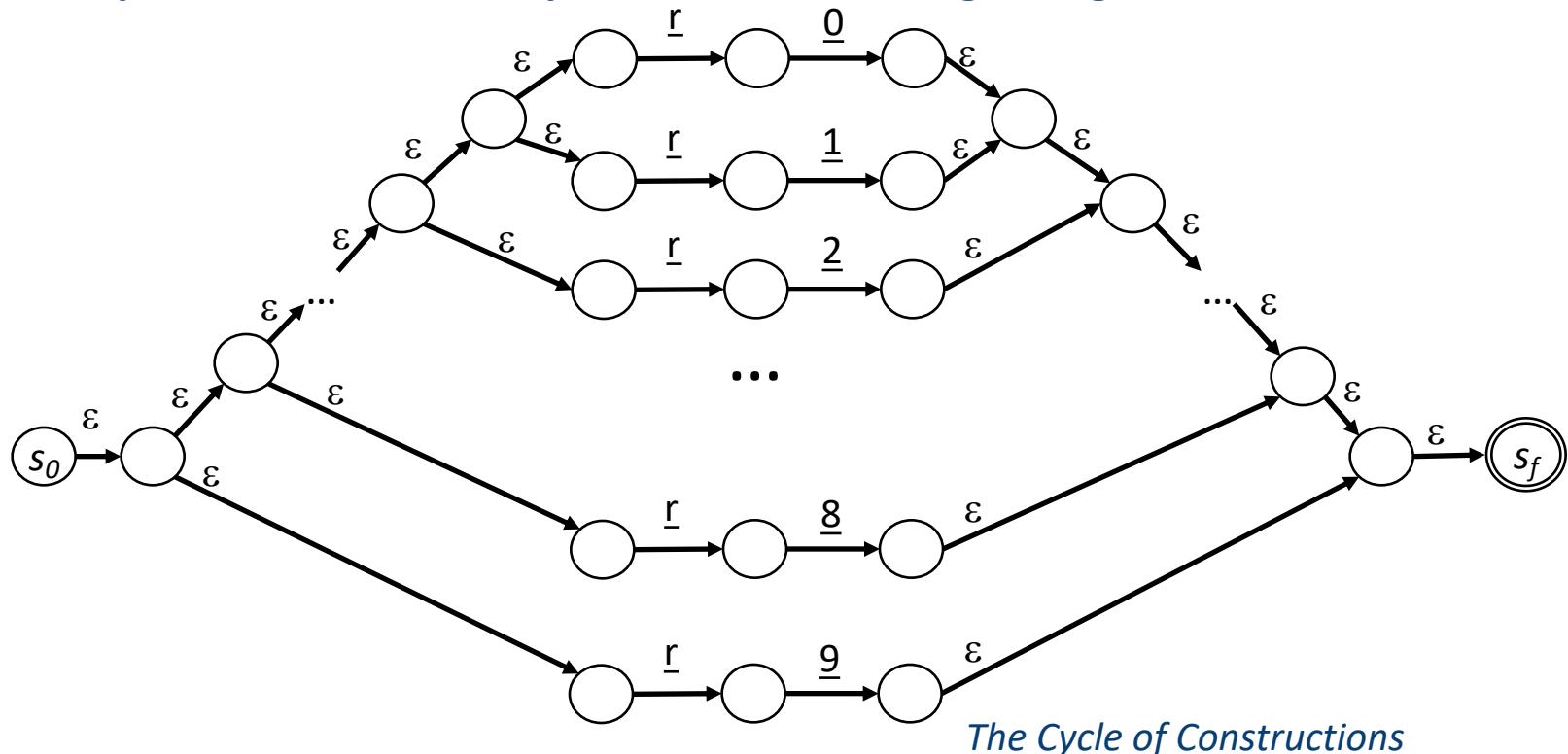
The Cycle of Constructions



# Abbreviated Register Specification



Thompson's construction produces something along these lines



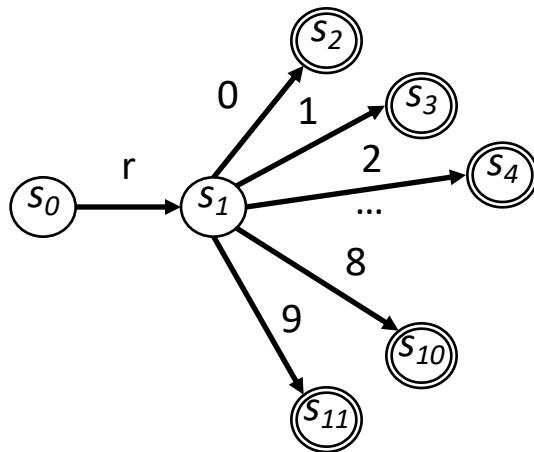
To make the example fit, we have eliminated some of the  $\epsilon$ -transitions, e.g., between r and 0



# Abbreviated Register Specification



Applying the subset construction yields



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

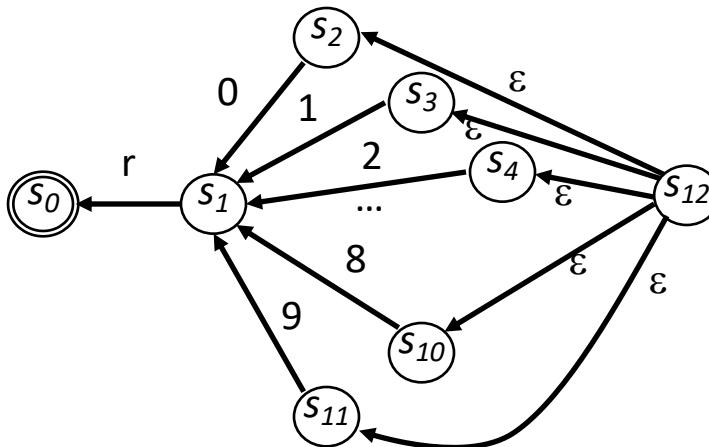
The Cycle of Constructions



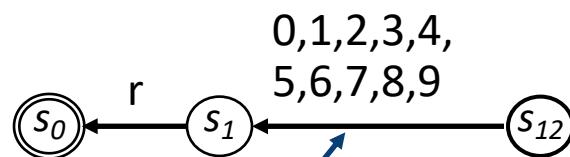
# Abbreviated Register Specification



## Applying Brzozowski's algorithm, step 1



Reversed NFA



After Subset Construction

The Cycle of Constructions

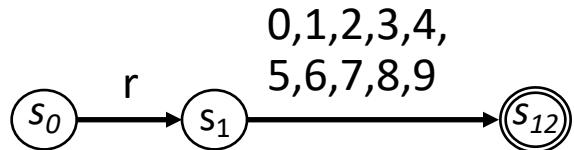
Technically, this edge shows up as 10 edges, which need to be combined...



# Abbreviated Register Specification



Brzozowski, step 2 reverses that DFA and subsets it again

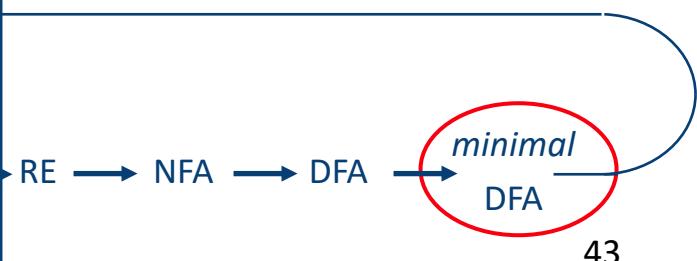


A skilled human might build this **DFA**

## The Critical Point:

- The construction will build a minimal **DFA**
- The size of the **DFA** relates to the language described by the **RE**, not the size of the **RE**
- The result is a **DFA**, so it has **O(1)** cost per character
- The compiler writer can use the “most natural” or “most intuitive” **RE**

## The Cycle of Constructions





# Where are we? Why are we doing this?

**RE → NFA** (*Thompson's construction*) ✓

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

**NFA → DFA** (*subset construction*) ✓

- Build the simulation

The Cycle of Constructions



**DFA → Minimal DFA**

- Hopcroft's algorithm
- Brzozowski's algorithm ✓

**DFA → RE**

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

## One Last Slide

# RE Back to DFA

The Wikipedia page on “Kleene’s algorithm” is pretty good. It also contains a link to Kleene’s 1956 paper. This form of the algorithm is usually attributed to McNaughton and Yamada in 1960.



## Kleene’s Construction

```
for i ← 0 to |D| - 1; // label each immediate path
    for j ← 0 to |D| - 1;
         $R^0_{ij} \leftarrow \{ a \mid \delta(d_i, a) = d_j \};$ 
        if (i = j) then
             $R^0_{ii} = R^0_{ii} \mid \{\epsilon\};$ 

for k ← 0 to |D| - 1; // label nontrivial paths
    for i ← 0 to |D| - 1;
        for j ← 0 to |D| - 1;
             $R^k_{ij} \leftarrow R^{k-1}_{ik} (R^{k-1}_{kk})^* R^{k-1}_{kj} \mid R^{k-1}_{ij}$ 

 $L \leftarrow \{ \} \quad // union labels of paths from$ 
For each final state  $s_i$  //  $s_0$  to a final state  $s_i$ 
     $L \leftarrow L \mid R^{|D|-1}_{0i}$ 
```

$R^k_{ij}$  is the set of paths from  $i$  to  $j$  that include no state higher than  $k$

Adaptation of all points, all paths,  
low cost algorithm

COMP 412, Fall 2017

The Cycle of Constructions

