



RICE

Early due date?
Friday or Saturday?

COMP 412
FALL 2013

Lexical Analysis — Part III
From NFA to DFA:
the Subset Construction

Comp 412

*With a quick look at Brzozowski's
Minimization Algorithm*

Copyright 2013, 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.



Where are we? Why are we doing this?

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

- Build an NFA for each term
- Combine them with ε -moves

NFA \rightarrow DFA (*subset construction*)

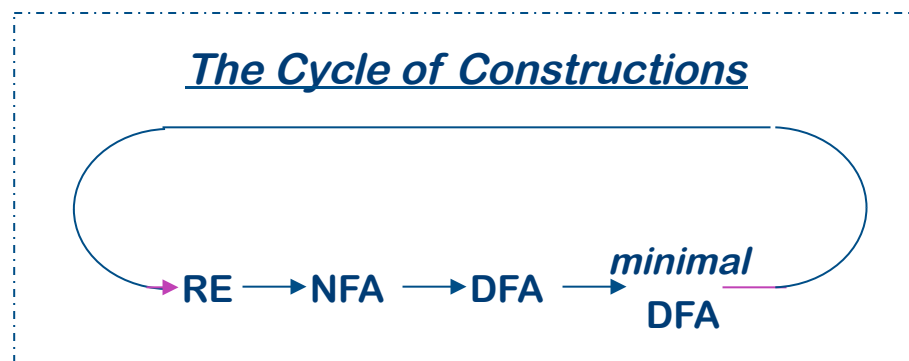
- Build the simulation

DFA \rightarrow Minimal DFA

- Hopcroft's algorithm

DFA \rightarrow RE

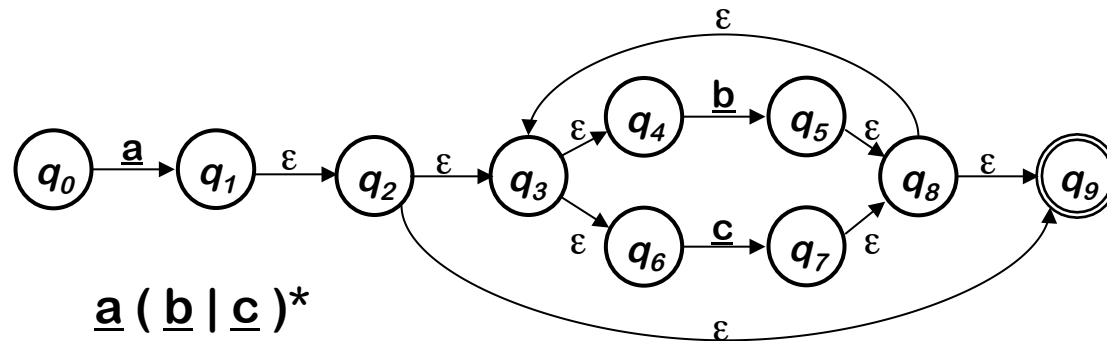
- All pairs, all paths problem
- Union together paths from s_0 to a final state



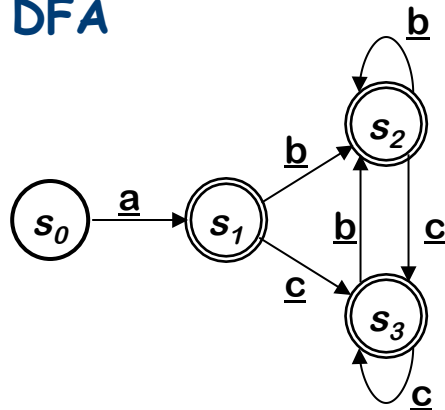


Simulating an NFA with a DFA

NFA



DFA



Where the mapping between NFA states and DFA states is:

DFA	NFA
s_0	q_0
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$
s_3	$q_7, q_8, q_9, q_3, q_4, q_6$



NFA \rightarrow DFA with Subset Construction

Need to build a simulation of the NFA

Two key functions

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

The algorithm:

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

Sounds more complex than it is...

Any DFA state containing a final state of the NFA becomes a final state of the DFA.



NFA \rightarrow DFA with Subset Construction

The algorithm:

```
 $s_0 \leftarrow \varepsilon\text{-closure}(\{n_0\})$   
 $S \leftarrow \{ s_0 \}$   
 $W \leftarrow \{ s_0 \}$   
while (  $W \neq \emptyset$  )  
    select and remove  $s$  from  $W$   
    for each  $\alpha \in \Sigma$   
         $t \leftarrow \varepsilon\text{-closure}(\text{Move}(s, \alpha))$   
         $T[s, \alpha] \leftarrow t$   
        if (  $t \notin S$  ) then  
            add  $t$  to  $S$   
            add  $t$  to  $W$ 
```

Let's think about why this works

s_0 is a set of states
 S & W are sets of sets of states

The algorithm halts:

1. S contains no duplicates (test before adding)
2. $2^{\{\text{NFA states}\}}$ 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

This test is a little tricky



NFA \rightarrow 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



$\underline{c})^*$:

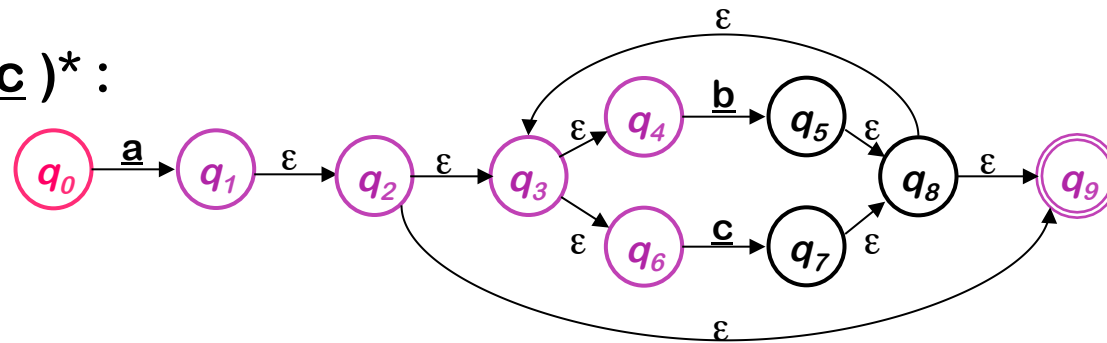
```
graph LR; q0((q0)) -- a --> q1((q1)); q1 -- ε --> q2((q2)); q2 -- ε --> q3((q3)); q3 -- ε --> q4((q4)); q3 -- ε --> q6((q6)); q4 -- b --> q5((q5)); q5 -- ε --> q8((q8)); q6 -- c --> q7((q7)); q7 -- ε --> q8; q8 -- ε --> q9(((q9))); q3 -- ε --> q8; q2 -- ε --> q9;
```

Comp 412, Fall 2013

NFA \rightarrow DFA with Subset Construction



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

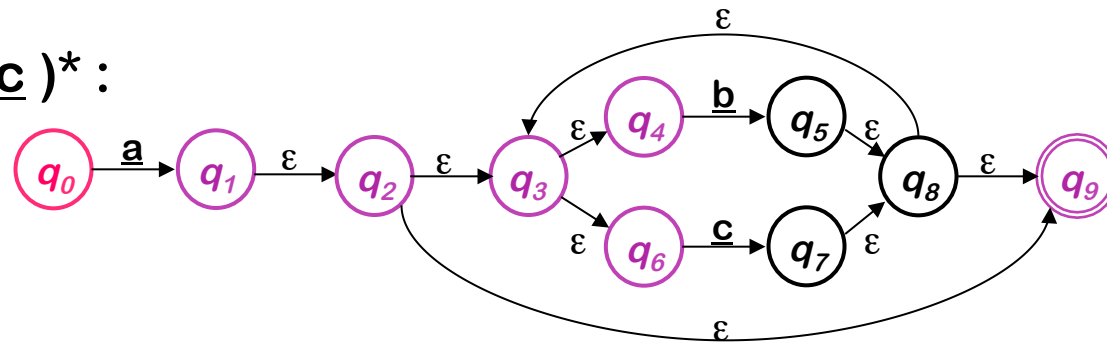


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$		

NFA \rightarrow DFA with Subset Construction



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

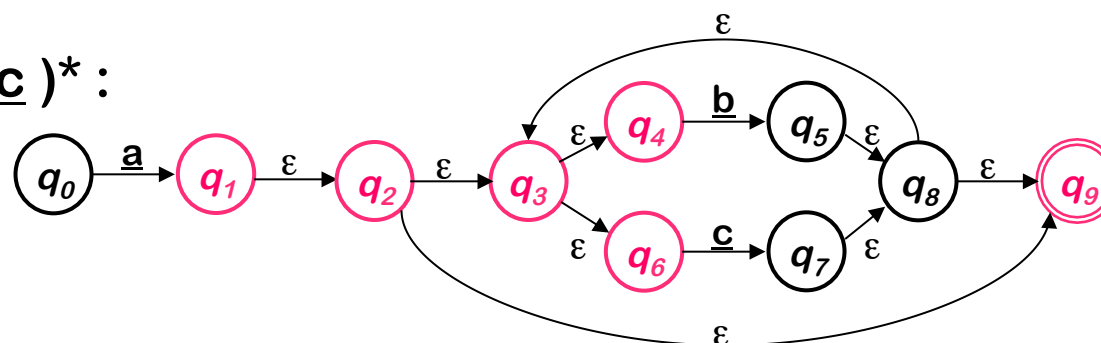


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3,$ q_4, q_6, q_9	none	none



NFA → DFA with Subset Construction

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

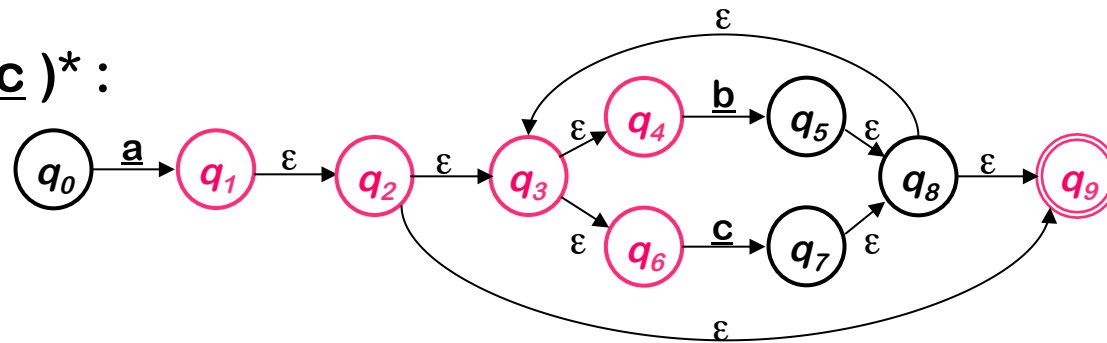


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$			

NFA \rightarrow DFA with Subset Construction



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

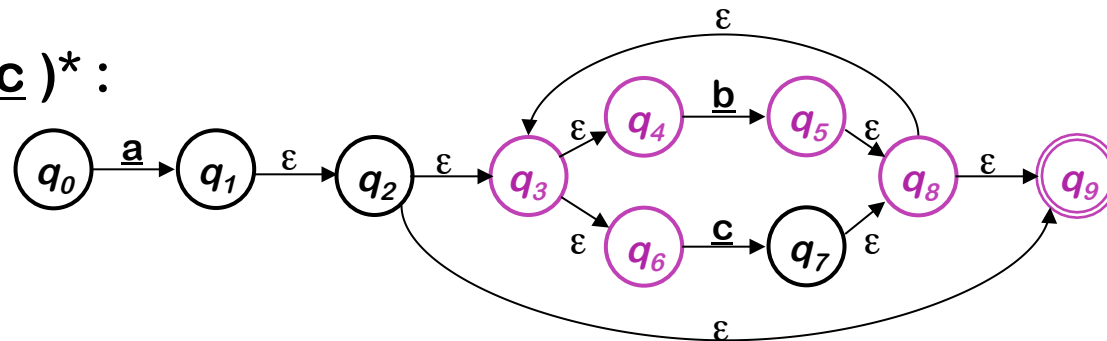


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none		

NFA \rightarrow DFA with Subset Construction



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

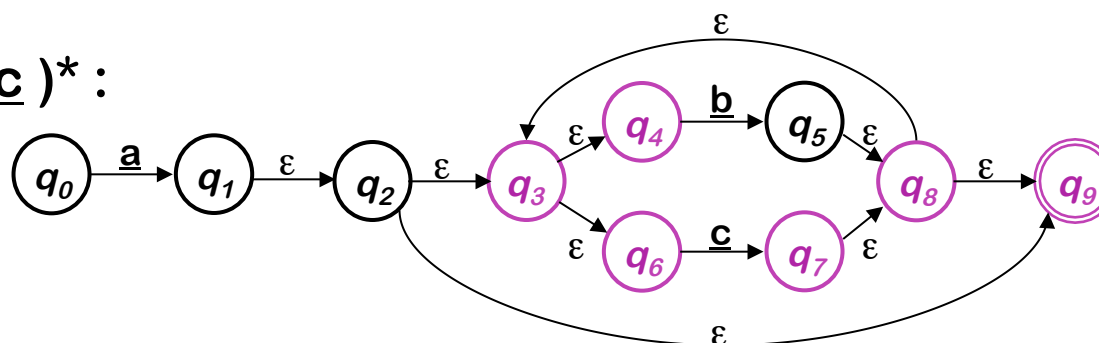


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	



NFA \rightarrow DFA with Subset Construction

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

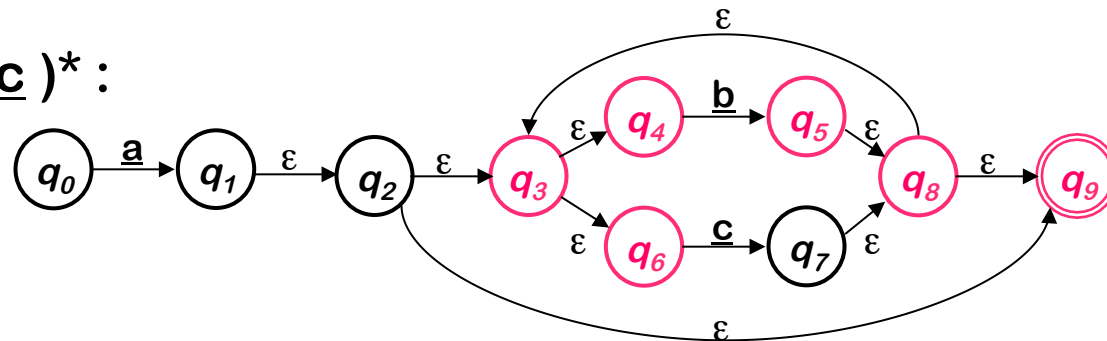


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$

NFA \rightarrow DFA with Subset Construction



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

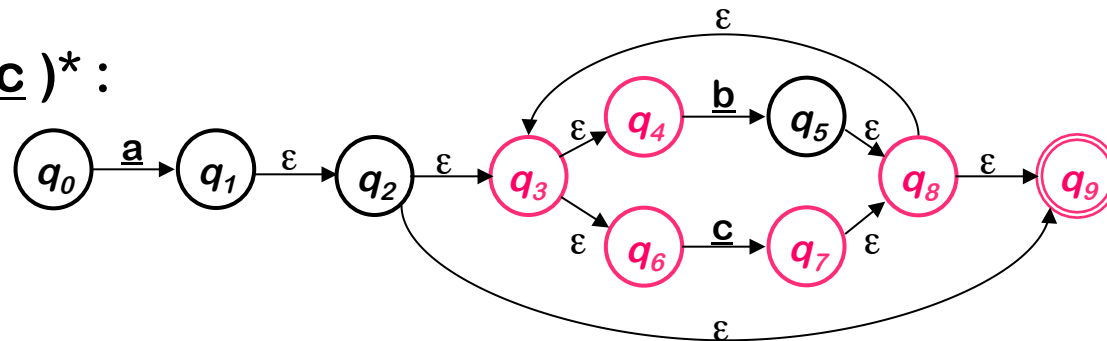


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$			

NFA → DFA with Subset Construction



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

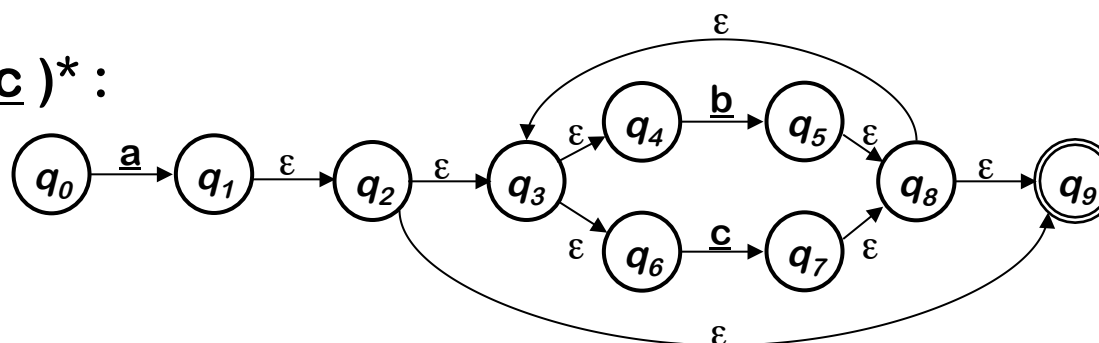


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$			
s_3	$q_7, q_8, q_9, q_3, q_4, q_6$			



NFA \rightarrow DFA with Subset Construction

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

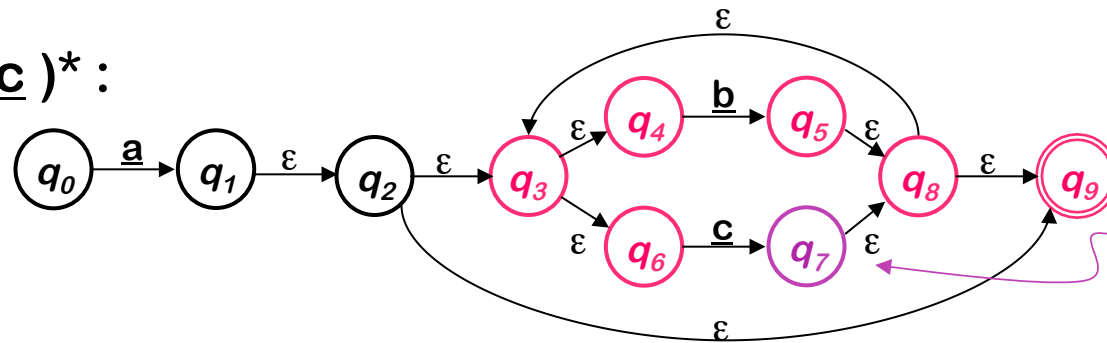


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$	none		
s_3	$q_7, q_8, q_9, q_3, q_4, q_6$	none		

NFA → DFA with Subset Construction



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



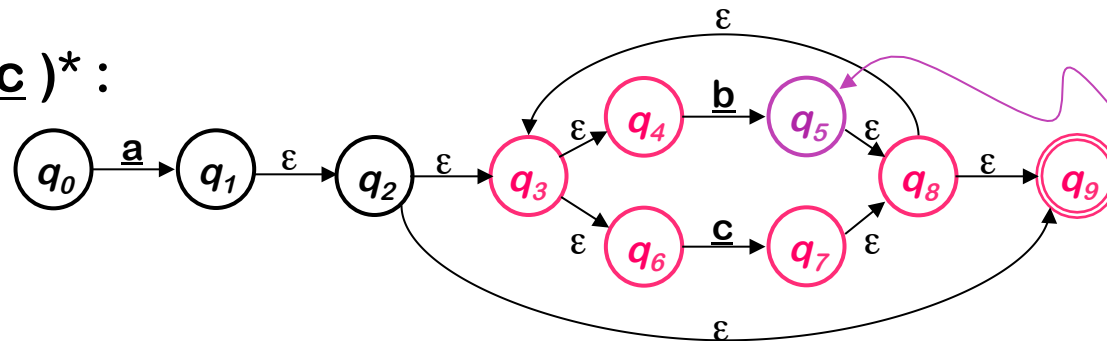
q_7 is the core state of s_3

States		ϵ -closure(Move($s, *$))		
DFA	NFA	<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		

NFA → DFA with Subset Construction



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



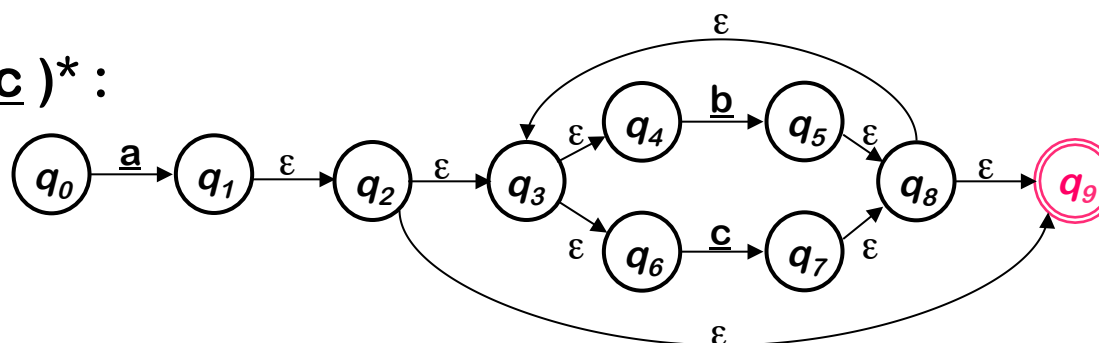
q₅ is the core state of s₂

States		ϵ -closure(Move(s,*))		
DFA	NFA	<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



NFA → DFA with Subset Construction

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

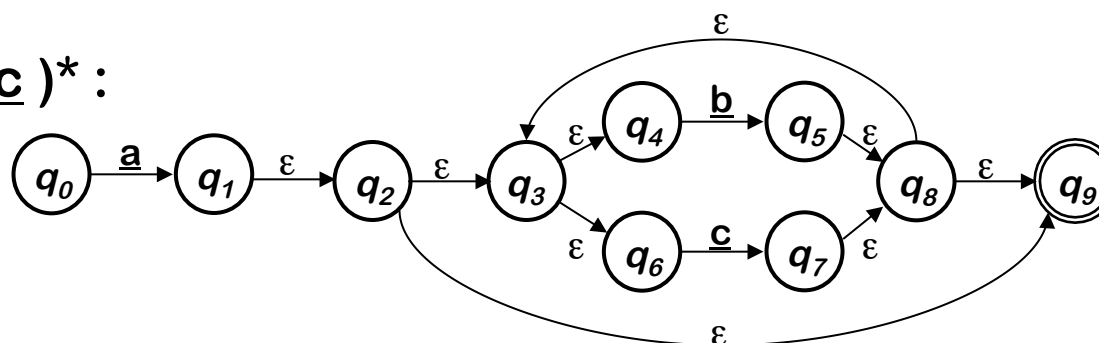


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_9$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_9$	none	$q_5, q_8, q_9, q_3, q_4, q_6$	$q_7, q_8, q_9, q_3, q_4, q_6$
s_2	$q_5, q_8, q_9, q_3, q_4, q_6$	none	s_2	s_3
s_3	$q_7, q_8, q_9, q_3, q_4, q_6$	none	s_2	s_3



NFA → DFA with Subset Construction

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

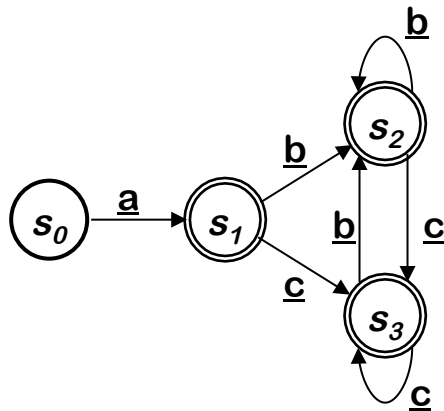


States		ϵ -closure(Move($s, *$))		
DFA	NFA	<u>a</u>	<u>b</u>	<u>c</u>
s_0	q_0	s_1	none	none
s_1	$q_1, q_2, q_3,$ q_4, q_6, q_9	none	s_2	s_3
s_2	$q_5, q_8, q_9,$ q_3, q_4, q_6	none	s_2	s_3
s_3	$q_7, q_8, q_9,$ q_3, q_4, q_6	none	s_2	s_3



NFA \rightarrow DFA with Subset Construction

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



	<u>a</u>	<u>b</u>	<u>c</u>
s_0	s_1	<i>none</i>	<i>none</i>
s_1	<i>none</i>	s_2	s_3
s_2	<i>none</i>	s_2	s_3
s_3	<i>none</i>	s_2	s_3

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





Where are we? Why are we doing this?

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

- Build an NFA for each term
- Combine them with ε -moves

NFA \rightarrow DFA (*subset construction*) ✓

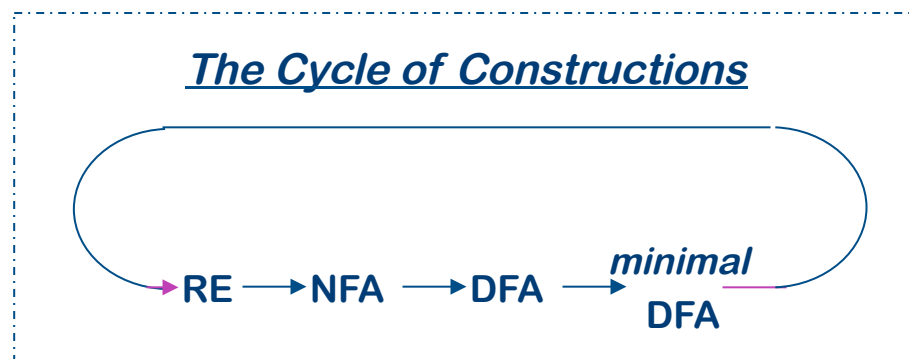
- Build the simulation

DFA \rightarrow Minimal DFA \Leftarrow

- Hopcroft's algorithm

DFA \rightarrow RE

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



Not enough time to teach Hopcroft's algorithm today

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

Proof: Assume first that $x = \sigma_0 \sigma_1 \dots \sigma_{n-1}$ is in $T(\mathfrak{A})$ and let s_0, s_1, \dots, s_n be a sequence of internal states of \mathfrak{A} such that s_0 is in S_0 and s_n is in F , 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, $s_{n-k} = s'_{k-1}$ 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.

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

Assume next that a tape $x = \sigma_0 \sigma_1 \dots \sigma_{n-1}$ is in $T(\mathfrak{D}(\mathfrak{A}))$. Let for each $k \leq n$, $t_k = N(t_{k-1}, \sigma_{k-1})$. 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

$$t_{n-1} = N(t_{n-2}, \sigma_{n-2}),$$

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 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 is in S_0 and s_n is in F , 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, $s_{n-k} = s'_{k-1}$ 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

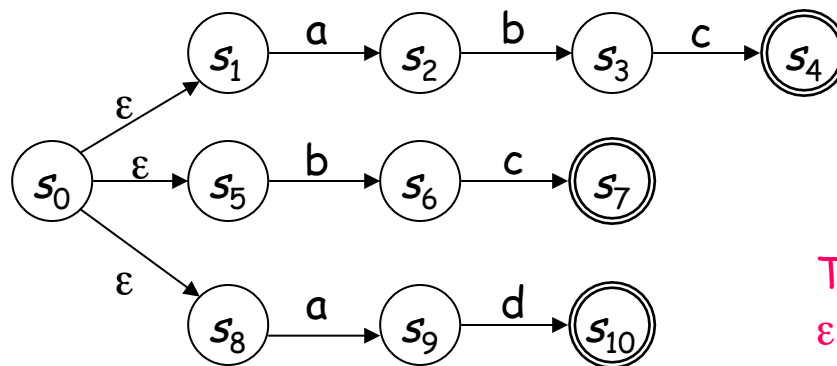
Simplifying a result due originally to Kleene, Myhill in unpublished work has shown that the class \mathcal{I} can be characterized as the least class of sets of tapes containing the finite sets and closed under some simple operations on sets of tapes. We indicate here a different proof using



Alternative Approach to DFA Minimization

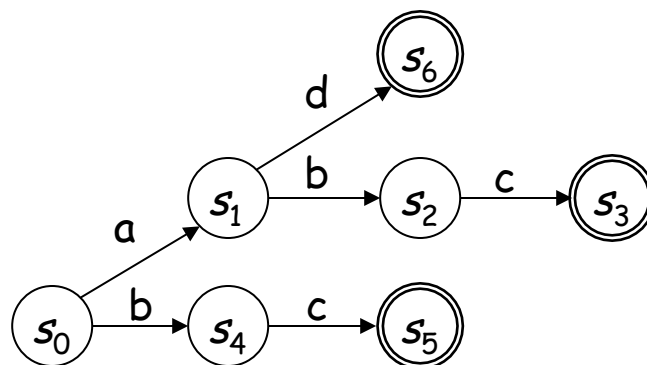
The Intuition

- The subset construction merges prefixes in the NFA



abc | bc | ad

Thompson's construction would leave ϵ -transitions between each single-character automaton



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



Alternative Approach to DFA Minimization

Idea: use the subset construction twice

- For an NFA N
 - Let $reverse(N)$ be the NFA constructed by making initial states final (& vice-versa) and reversing the edges
 - Let $subset(N)$ be the DFA that results from applying the subset construction to N
 - Let $reachable(N)$ be N after removing all states that are not reachable from the initial state
- Then,

$reachable(subset(reverse[reachable(subset(reverse(N))])))$

is the minimal DFA that implements N [Brzozowski, 1962]

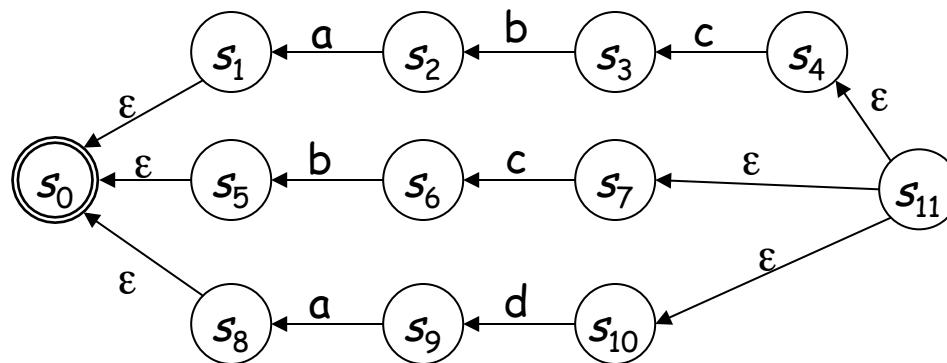
*This result is not intuitive, but it is true.
Neither algorithm dominates the other.*



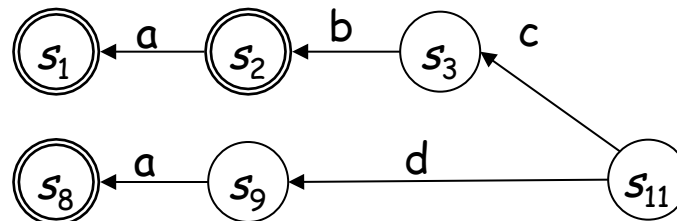
Alternative Approach to DFA Minimization

Step 1

- The subset construction on $reverse(NFA)$ merges suffixes in original NFA



Reversed NFA



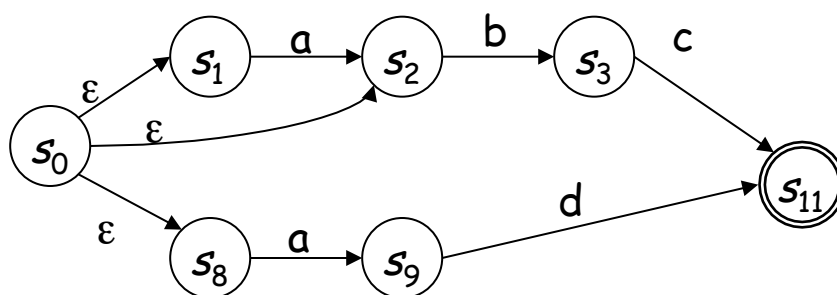
subset(reverse(NFA))



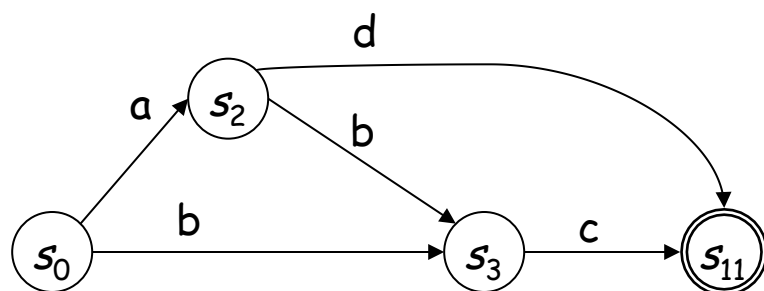
Alternative Approach to DFA Minimization

Step 2

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



Reverse it, again



Minimal DFA

And subset it, again

The Cycle of Constructions

