



RICE

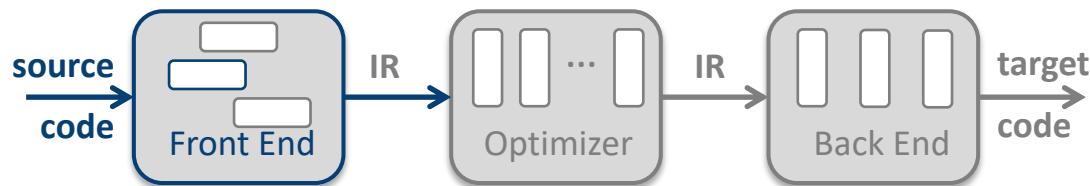
Ignore § 2.4.4 in EaC2e.  
I will post a replacement  
section on the course web site.

COMP 412  
FALL 2018

## Lexical Analysis, III

Comp 412

Minor corrections  
applied after class



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



# 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

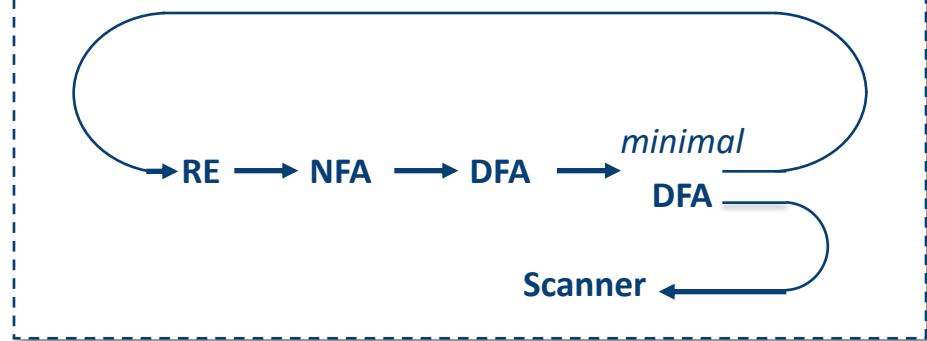
**Minimal DFA → Scanner**

- 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*





# DFA Minimization

## The Big Picture

- Discover sets of behaviorally equivalent states in the DFA
- Represent each such set with a single new state

Two states  $s_i$  and  $s_j$  are **behaviorally equivalent if and only if**:

*Recursive definition*

- $\forall c \in \Sigma$ , transitions from  $s_i$  &  $s_j$  on  $c$  lead to equivalent states
- The set of paths leading from  $s_i$  &  $s_j$  are equivalent

A **partition**  $P$  of a set  $S$ :

- A collection of subsets of  $P$  such that each state  $s$  is in exactly one  $p_i \in P$

The algorithm iteratively constructs partitions of the DFA's set of states

We want a partition  $P = \{ p_0, p_1, p_2, \dots p_n \}$  of  $D$  that has two properties:

1. If  $d_i$  &  $d_j \in p_s$  and  $c$  takes  $d_i \rightarrow d_x$  and  $d_j \rightarrow d_y$ , then  $d_x$  &  $d_y \in p_t$ ,  $\forall c, i, j, s, t$
2. If  $d_i$  &  $d_j \in p_s$  and  $d_i \in F$  then  $d_j \in F$

$D$  is the set of states for the DFA:  $(D, \Sigma, \delta, s_0, D_A)$

# DFA Minimization

Maximally sized sets  $\Rightarrow$   
minimal number of sets



## Details of the algorithm

- Group states into maximally-sized initial sets, *optimistically* (property 2)
- Iteratively subdivide those sets, based on transition graph (property 1)
- States that remain grouped together are equivalent

Initial partition:  $P_0$  has two sets:  $\{D_A\}$  &  $\{D - D_A\}$

$$D = (D, \Sigma, \delta, s_0, D_A)$$

| final states | other states |
|--------------|--------------|
|--------------|--------------|

## Property 1 provides the basis for refining, or splitting, the sets

- Assume  $s_i$  &  $s_j \in p_s$ , and  $\delta(s_i, \underline{a}) = s_x$ , &  $\delta(s_j, \underline{a}) = s_y$
- If  $s_x$  &  $s_y$  are not in the same set  $p_t$ , then  $p_s$  must be split
  - COROLLARY:  $s_i$  has transition on  $\underline{a}$ ,  $s_j$  does not  $\Rightarrow \underline{a}$  splits  $p_s$
- A single state in a DFA cannot have two transitions on  $\underline{a}$ 
  - Each  $p_s$  will become a DFA state

Algorithm actually works backward; it looks at what transitions enter  $p$  on character  $c$ , and uses that to split the partition  $q$  where those edges begin.



# DFA Minimization Algorithm (Worklist version)

```

Worklist  $\leftarrow \{D_A, \{D - D_A\}\}$ 
Partition  $\leftarrow \{D_A, \{D - D_A\}\}$ 
While (Worklist  $\neq \emptyset$ ) do
    select a set S from Worklist and remove it
    for each  $\alpha \in \Sigma$  do
        Image  $\leftarrow \{x \mid \delta(x, \alpha) \in S\}$ 
        for each  $q \in \text{Partition}$  that has a state in Image do
             $q_1 \leftarrow q \cap \text{Image}$ 
             $q_2 \leftarrow q - q_1$ 
            if  $q_2 \neq \emptyset$  then
                remove q from Partition
                Partition  $\leftarrow \text{Partition} \cup q_1 \cup q_2$ 
                if  $q \in \text{Worklist}$  then
                    remove q from Worklist
                    Worklist  $\leftarrow \text{Worklist} \cup q_1 \cup q_2$ 
                else if  $|q_1| \leq |q_2|$ 
                    then Worklist  $\leftarrow \text{Worklist} \cup q_1$ 
                    else Worklist  $\leftarrow \text{Worklist} \cup q_2$ 
                if  $s = q$  then
                    break; // cannot keep working on s
            end if
        end for
    end for
end while

```

Image is the set of states that have a transition into S on  $\alpha$ :  $\delta^{-1}(S, \alpha)$

$q_1$  is the subset of  $q$  that transitions to S on  $\alpha$   
 $q_2$  is the rest of  $q$

And, as an implementation hint, if we just split S — that is, S was  $q$  & it split — we need a new S

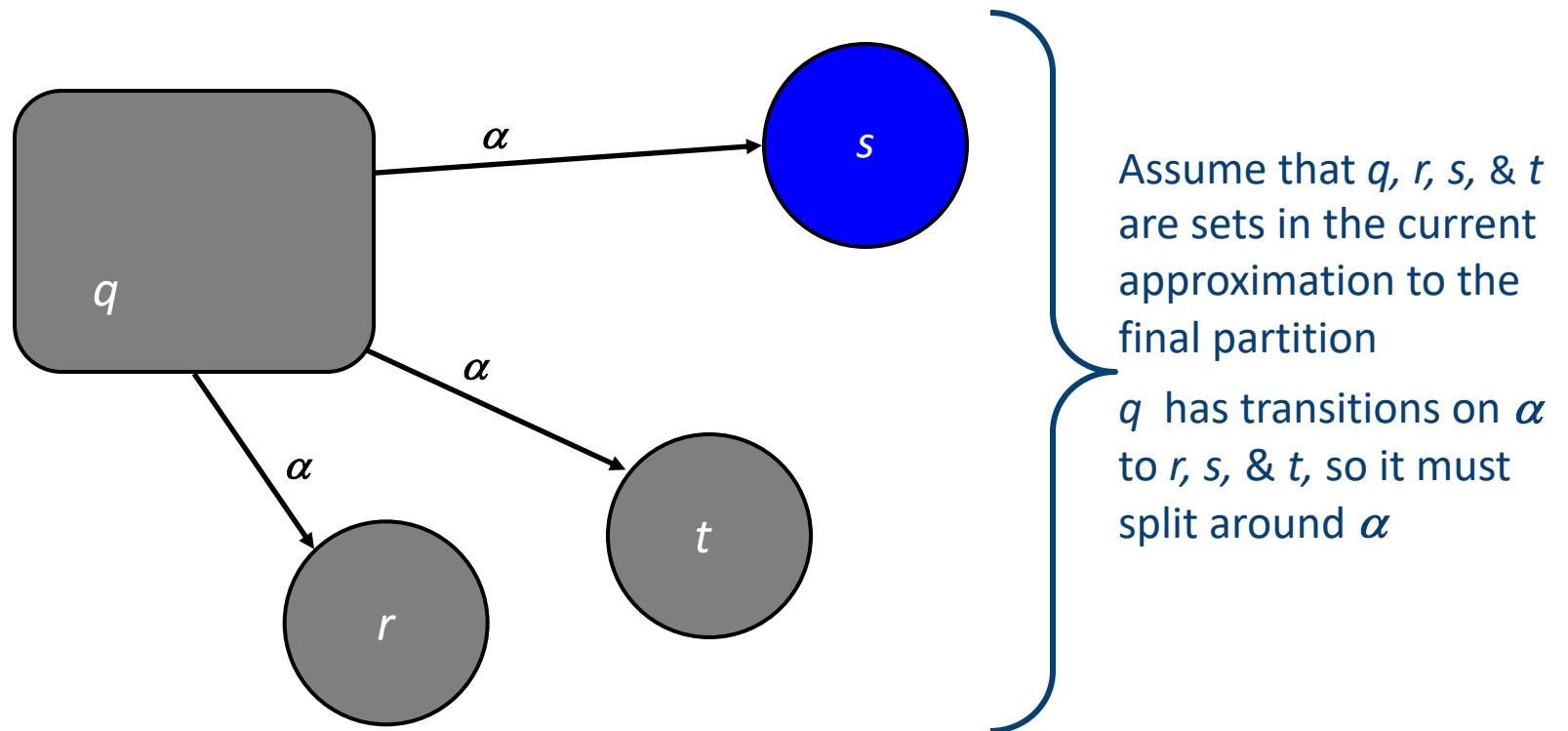
“split q”

adjust Worklist

# Key Idea: Splitting Q Around Transitions on $\alpha$



## Partitioning Q around S

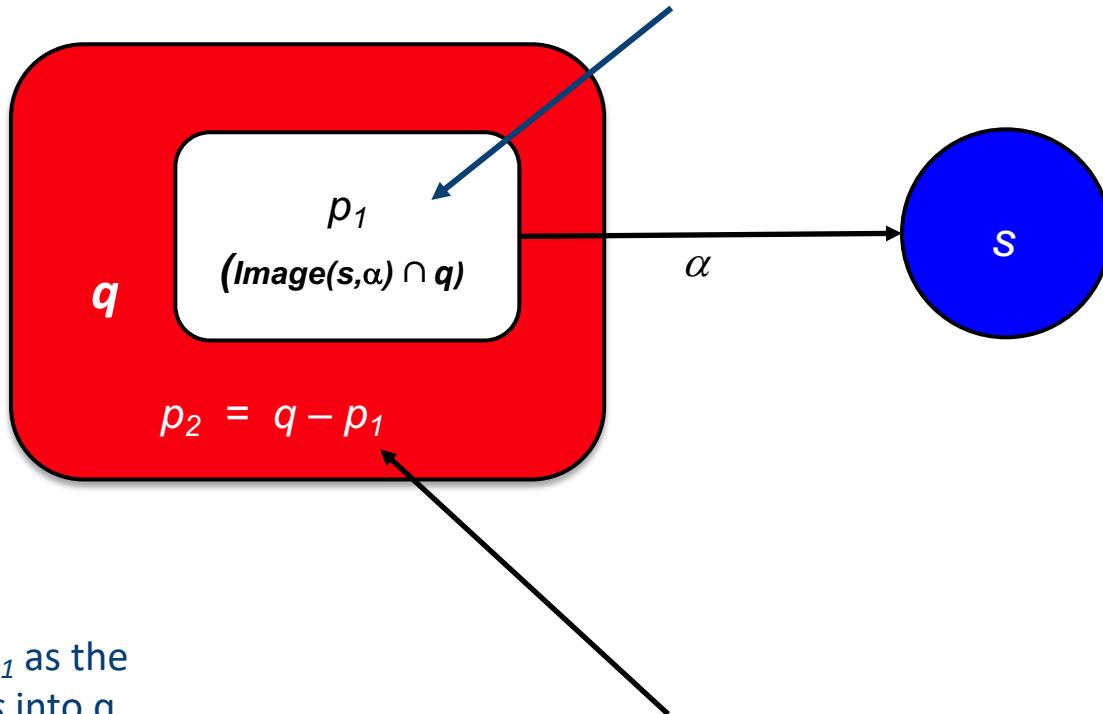


*As the algorithm considers  $s$  and  $\alpha$ , it will split  $q$ .*



# Key Idea: Splitting $q$ around $s$ and $\alpha$

Find maximal subset of  $q$  ( $p_1$ ) that has an  $\alpha$ -transition into  $s$



Think of  $p_1$  as the image of  $s$  into  $q$  under the inverse of the transition function:

$$p_1 \leftarrow \delta^{-1}(s, \alpha) \cap q$$

$p_2$  must have an  $\alpha$ -transition to one or more other states in one or more other partitions (e.g.,  $r$  &  $s$ ), or states with no  $\alpha$ -transitions.

Otherwise,  $q$  does not split!



# DFA Minimization Algorithm (Worklist version)

```

Worklist  $\leftarrow \{D_A, \{D - D_A\}\}$ 
Partition  $\leftarrow \{D_A, \{D - D_A\}\}$ 
While (Worklist  $\neq \emptyset$ ) do
    select a set S from Worklist and remove it
    for each  $\alpha \in \Sigma$  do
        Image  $\leftarrow \{x \mid \delta(x, \alpha) \in S\}$ 
        for each  $q \in \text{Partition}$  that has a state in Image do
             $q_1 \leftarrow q \cap \text{Image}$ 
             $q_2 \leftarrow q - q_1$ 
            if  $q_2 \neq \emptyset$  then
                remove q from Partition
                Partition  $\leftarrow \text{Partition} \cup q_1 \cup q_2$ 
                if  $q \in \text{Worklist}$  then
                    remove q from Worklist
                    Worklist  $\leftarrow \text{Worklist} \cup q_1 \cup q_2$ 
                else if  $|q_1| \leq |q_2|$ 
                    then Worklist  $\leftarrow \text{Worklist} \cup q_1$ 
                    else Worklist  $\leftarrow \text{Worklist} \cup q_2$ 
                if  $s = q$  then
                    break; // cannot keep working on s
            end if
        end for
    end for
end while

```

Projection is the set of states that have a transition into  $S$  on  $\alpha$ :  
 $\delta^{-1}(S, \alpha)$

$p_1$  is the subset of  $q$  that transitions to  $S$  on  $\alpha$   
 $p_2$  is the rest of  $q$

And, as an implementation hint, if we just split  $S$  — that is,  $S$  was  $q$  & it split — we need a new  $S$



# DFA Minimization Algorithm (Worklist version)

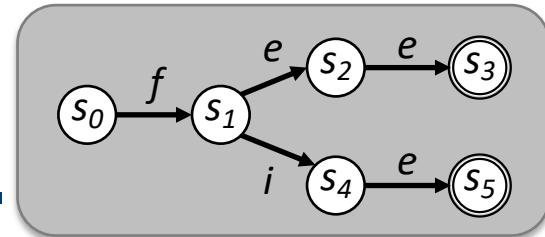
```
Worklist  $\leftarrow \{D_A, \{D - D_A\}\}$   
Partition  $\leftarrow \{D_A, \{D - D_A\}\}$ 
```

One last hack ...

```
While (Worklist  $\neq \emptyset$ ) do  
    select a set S from Worklist and remove it  
    for each  $\alpha \in \Sigma$  do  
        Image  $\leftarrow \{x \mid \delta(x, \alpha) \in S\}$   
        for each  $q \in \text{Partition}$  that has a state in Image do  
             $q_1 \leftarrow q \cap \text{Image}$   
             $q_2 \leftarrow q - q_1$   
            if  $q_2 \neq \emptyset$  then  
                remove  $q$  from Partition  
                Partition  $\leftarrow \text{Partition} \cup q_1 \cup q_2$   
                if  $q \in \text{Worklist}$  then  
                    remove  $q$  from Worklist  
                    Worklist  $\leftarrow \text{Worklist} \cup q_1 \cup q_2$   
                else if  $|q_1| \leq |q_2|$   
                    then Worklist  $\leftarrow \text{Worklist} \cup q_1$   
                    else Worklist  $\leftarrow \text{Worklist} \cup q_2$   
                if  $s = q$  then  
                    break; // cannot keep working on s
```

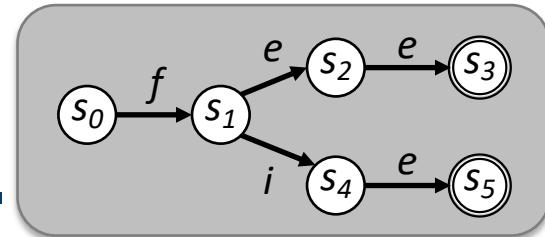
If  $q$  is a singleton, we can skip the body of the loop because a singleton cannot split.

# A Detailed Example



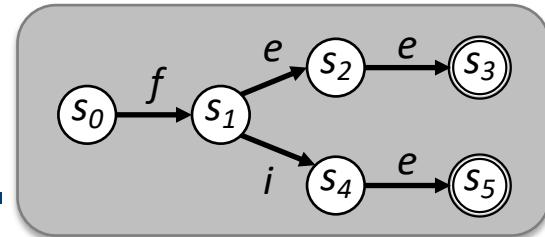
| Step | Partition                        | W'list     | s | c | Image | q | q <sub>1</sub> | q <sub>2</sub> | Action |
|------|----------------------------------|------------|---|---|-------|---|----------------|----------------|--------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$ | $p_0, p_1$ | — | — | —     | — | —              | —              | —      |

# A Detailed Example



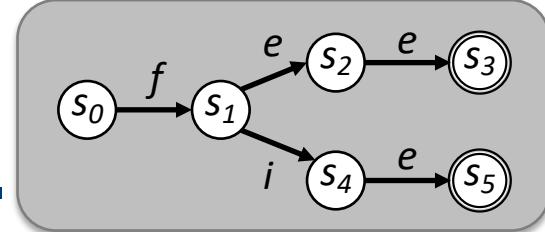
| Step | Partition                        | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                     |
|------|----------------------------------|------------|-------|---|------------|-------|----------------|----------------|----------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$ | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                          |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$ | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <i>split p<sub>1</sub></i> |

# A Detailed Example



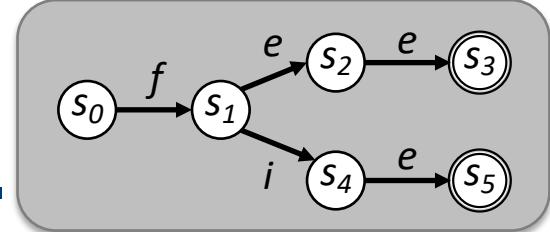
| Step | Partition                                  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                     |
|------|--|------------|-------|---|------------|-------|----------------|----------------|----------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$           | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                          |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$           | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <i>split p<sub>1</sub></i> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$ | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | <i>none</i>                |
|      |  | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | <i>none</i>                |

# A Detailed Example



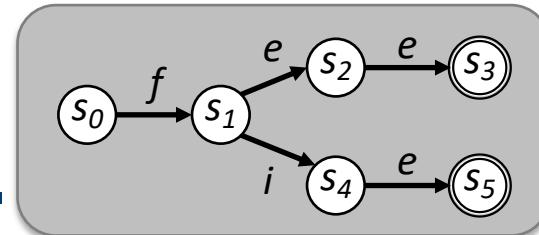
| Step | Partition                                  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                     |
|------|--|------------|-------|---|------------|-------|----------------|----------------|----------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$           | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                          |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$           | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <i>split p<sub>1</sub></i> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$ | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | <i>none</i>                |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$ | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | <i>none</i>                |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$ | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <i>split p<sub>3</sub></i> |

# A Detailed Example



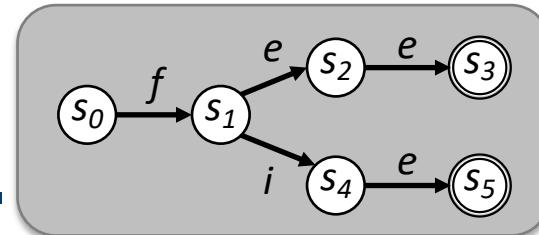
| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                        |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none                          |

# A Detailed Example



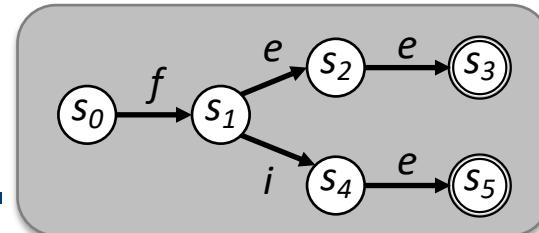
| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                        |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none                          |
| 3    |  | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none                          |

# A Detailed Example



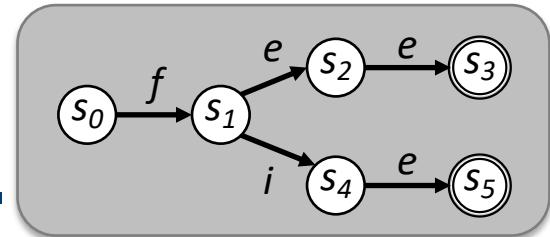
| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                        |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none                          |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | f | $s_0$      | $p_5$ | $s_0$          | ∅              | none                          |

# A Detailed Example



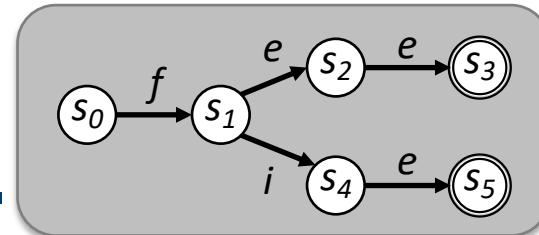
| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                        |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none                          |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | f | $s_0$      | $p_5$ | $s_0$          | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | i | ∅          | ∅     | ∅              | ∅              | none                          |

# A Detailed Example



| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action                        |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none                          |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | f | $s_0$      | $p_5$ | $s_0$          | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | i | ∅          | ∅     | ∅              | ∅              | none                          |
| 4    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | ∅          | $p_5$ | e | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | ∅          | $p_5$ | f | ∅          | ∅     | ∅              | ∅              | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | ∅          | $p_5$ | i | ∅          | ∅     | ∅              | ∅              | none                          |

# A Detailed Example

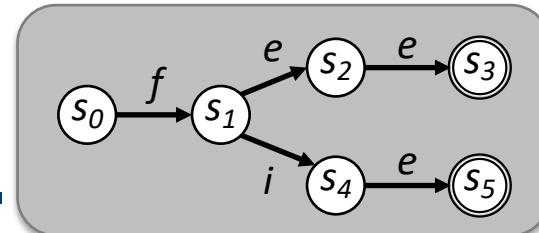


| Step | Partition  | W'list     | $s$   | $c$ | Image       | $q$         | $q_1$       | $q_2$       | Action                        |
|------|--|------------|-------|-----|-------------|-------------|-------------|-------------|-------------------------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | —   | —           | —           | —           | —           | —                             |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | $e$ | $s_2, s_4$  | $p_1$       | $s_2, s_4$  | $s_0, s_1$  | <b>split <math>p_1</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | $f$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | $i$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | none                          |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | $e$ | $s_1$       | $p_3$       | $s_1$       | $s_0$       | <b>split <math>p_3</math></b> |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | $f$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | $i$ | $s_1$       | $p_4$       | $s_1$       | $\emptyset$ | none                          |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | $e$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | $f$ | $s_0$       | $p_5$       | $s_0$       | $\emptyset$ | none                          |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ |            |       |     |             | $\emptyset$ | $\emptyset$ | $\emptyset$ | none                          |

## Reconstructing the DFA

- Each set in Partition forms a state
- For each line in the table where both  $q_1$  and  $s$ , have values, add an edge from  $q_1$  to  $s$  labelled  $c$

# A Detailed Example

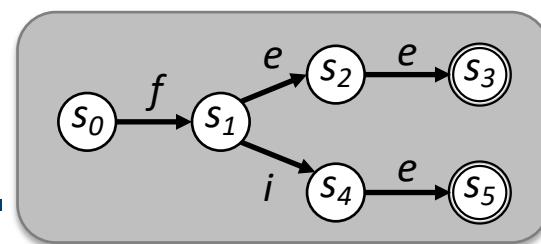
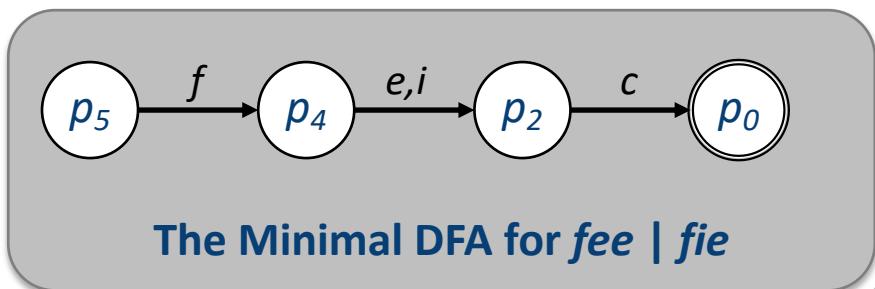


| Step | Partition  | W'list     | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action      |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —           |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | split $p_1$ |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none        |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | split $p_3$ |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none        |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none        |
|      |  | $p_5$      | $p_4$ | f | $s_0$      | $p_5$ | $s_0$          | ∅              | none        |

## Reconstructing the DFA

- Each set in Partition forms a state
- For each line in the table where  $q_1 \neq \emptyset$ , add an edge from  $q_1$  to  $s$  labelled  $c$

|   |   |   |   |   |      |
|---|---|---|---|---|------|
| i | ∅ | ∅ | ∅ | ∅ | none |
| e | ∅ | ∅ | ∅ | ∅ | none |
| f | ∅ | ∅ | ∅ | ∅ | none |
| i | ∅ | ∅ | ∅ | ∅ | none |



The Minimal DFA for *fee* | *fie*

| Step | Partition  | vv list    | s     | c | Image      | q     | q <sub>1</sub> | q <sub>2</sub> | Action      |
|------|--|------------|-------|---|------------|-------|----------------|----------------|-------------|
| 0    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_0, p_1$ | —     | — | —          | —     | —              | —              | —           |
| 1    | $p_0: \{3,5\}, p_1: \{0,1,2,4\}$                     | $p_1$      | $p_0$ | e | $s_2, s_4$ | $p_1$ | $s_2, s_4$     | $s_0, s_1$     | split $p_1$ |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_2, p_3$ | $p_0$ | f | ∅          | ∅     | ∅              | ∅              | none        |
|      |  | $p_2, p_3$ | $p_0$ | i | ∅          | ∅     | ∅              | ∅              | none        |
| 2    | $p_0: \{3,5\}, p_2: \{2,4\}, p_3: \{0,1\}$           | $p_3$      | $p_2$ | e | $s_1$      | $p_3$ | $s_1$          | $s_0$          | split $p_3$ |
|      | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_4, p_5$ | $p_2$ | f | ∅          | ∅     | ∅              | ∅              | none        |
|      |  | $p_4, p_5$ | $p_2$ | i | $s_1$      | $p_4$ | $s_1$          | ∅              | none        |
| 3    | $p_0: \{3,5\}, p_2: \{2,4\}, p_4: \{1\}, p_5: \{0\}$ | $p_5$      | $p_4$ | e | ∅          | ∅     | ∅              | ∅              | none        |
|      |  | $p_5$      | $p_4$ | f | $s_0$      | $p_5$ | $s_0$          | ∅              | none        |

### Reconstructing the DFA

- Each set in Partition forms a state
- For each line in the table where  $q_1 \neq \emptyset$ , add an edge from  $q_1$  to  $s$  labelled  $c$

|   |   |   |   |   |      |
|---|---|---|---|---|------|
| i | ∅ | ∅ | ∅ | ∅ | none |
| e | ∅ | ∅ | ∅ | ∅ | none |
| f | ∅ | ∅ | ∅ | ∅ | none |
| i | ∅ | ∅ | ∅ | ∅ | none |



# DFA Minimization Algorithm (Worklist version)

```
Worklist  $\leftarrow \{D_A, \{D - D_A\}\}$ 
Partition  $\leftarrow \{D_A, \{D - D_A\}\}$ 
While (Worklist  $\neq \emptyset$ ) do
    select a set S from Worklist and remove it
    for each  $\alpha \in \Sigma$  do
        Image  $\leftarrow \{x \mid \delta(x, \alpha) \in S\}$ 
        for each  $q \in \text{Partition}$  that has a state in Image
             $q_1 \leftarrow q \cap \text{Image}$ 
             $q_2 \leftarrow q - q_1$ 
            if  $q_2 \neq \emptyset$  then
                remove  $q$  from Partition
                Partition  $\leftarrow \text{Partition} \cup q_1 \cup q_2$ 
                if  $q \in \text{Worklist}$  then
                    remove  $q$  from Worklist
                    Worklist  $\leftarrow \text{Worklist} \cup q_1 \cup q_2$ 
                else if  $|q_1| \leq |q_2|$ 
                    then Worklist  $\leftarrow \text{Worklist} \cup q_1$ 
                    else Worklist  $\leftarrow \text{Worklist} \cup q_2$ 
                if  $s = q$  then
                    break; // cannot keep working on s
```

## Why does this algorithm halt?

- Fixed-point algorithm
- DFA has finite number of states
- Start with 2 sets in Partition
- Splitting breaks 1 set into 2 smaller ones but never makes a set larger
  - Monotone behavior
- Simple, finite limit on  $|\text{Partition}|$ ; it cannot be  $> |\text{States}|$
- Finite # steps, monotone increasing construction  $\Rightarrow$  algorithm halts



# DFA Minimization Algorithm (Worklist version)

```
Worklist  $\leftarrow \{D_A, \{D - D_A\}\}$ 
Partition  $\leftarrow \{D_A, \{D - D_A\}\}$ 
While (Worklist  $\neq \emptyset$ ) do
    select a set S from Worklist and remove it
    for each  $\alpha \in \Sigma$  do
        Image  $\leftarrow \{x \mid \delta(x, \alpha) \in S\}$ 
        for each  $q \in \text{Partition}$  that has a state in Image do
             $q_1 \leftarrow q \cap \text{Image}$ 
             $q_2 \leftarrow q - q_1$ 
            if  $q_2 \neq \emptyset$  then
                remove  $q$  from Partition
                Partition  $\leftarrow \text{Partition} \cup q_1 \cup q_2$ 
                if  $q \in \text{Worklist}$  then
                    remove  $q$  from Worklist
                    Worklist  $\leftarrow \text{Worklist} \cup q_1 \cup q_2$ 
                else if  $|q_1| \leq |q_2|$ 
                    then Worklist  $\leftarrow \text{Worklist} \cup q_1$ 
                    else Worklist  $\leftarrow \text{Worklist} \cup q_2$ 
                if  $s = q$  then
                    break; // cannot keep working on s
```

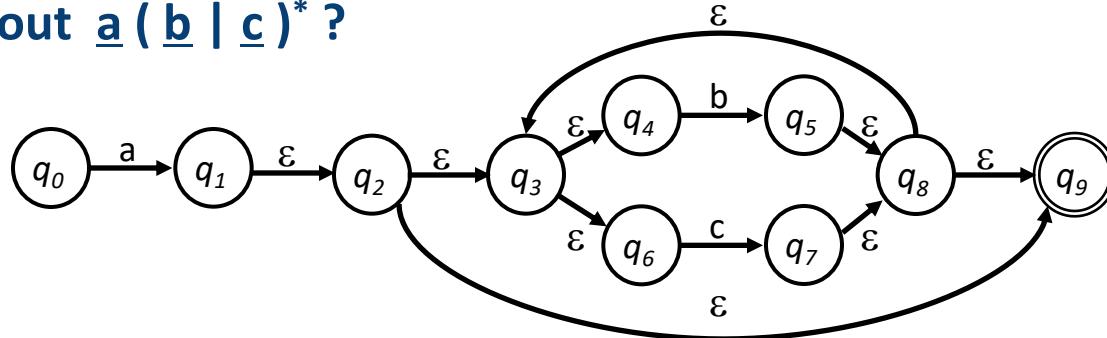
## One last hack ...

To make an implementation faster, it should maintain an efficient way to determine, for a given state, which set currently contain that state.



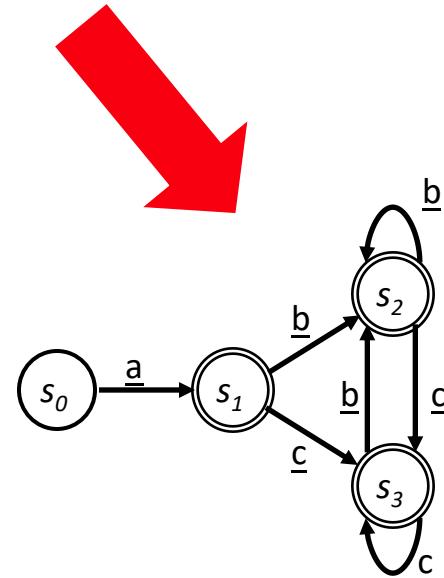
# DFA Minimization

What about  $\underline{a} (\underline{b} \mid \underline{c})^*$ ?



From the subset construction:

| States |                                | $\epsilon$ -closure( $\text{Move}(s, *)$ ) |          |          |
|--------|--------------------------------|--|----------|----------|
| 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                                       | $s_2$    | $s_3$    |

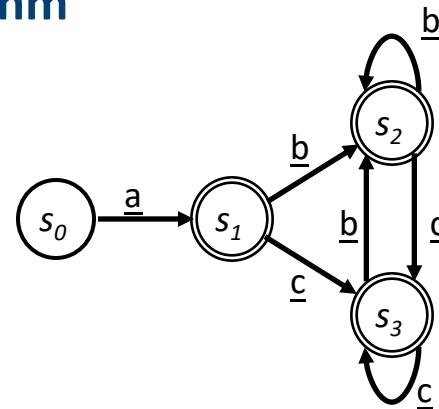


# DFA Minimization



## Applying Hopcroft's DFA minimization algorithm

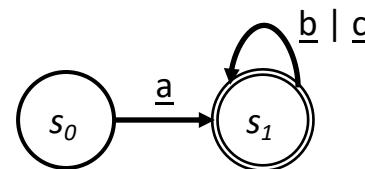
|                   |                                   | Split on |          |          |
|-------------------|-----------------------------------|----------|----------|----------|
| Current Partition |                                   | <u>a</u> | <u>b</u> | <u>c</u> |
| $P_0$             | $\{s_1,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



It splits no states after the initial partition

⇒ The minimal DFA has two states

- One for  $\{s_0\}$
- One for  $\{s_1, s_2, s_3\}$



Earlier, I suggested that a human would design a simpler automaton than Thompson's construction & the subset construction did.

Minimizing that DFA produces exactly the DFA that I claimed a human would design!

# Abbreviated Register Specification



Start with a regular expression

$r0 \mid r1 \mid r2 \mid r3 \mid r4 \mid r5 \mid r6 \mid r7 \mid r8 \mid r9$

Register names from zero to nine

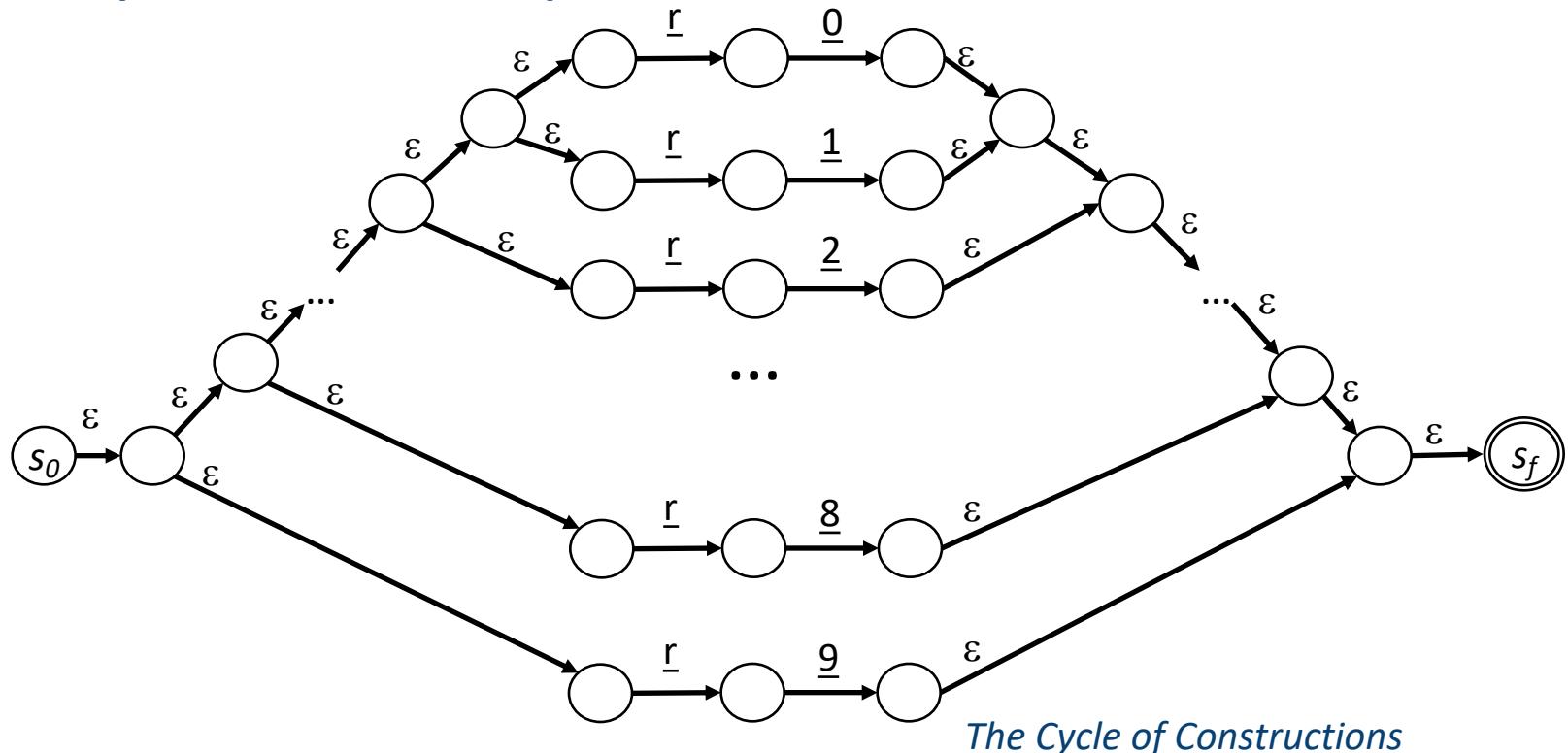
The Cycle of Constructions



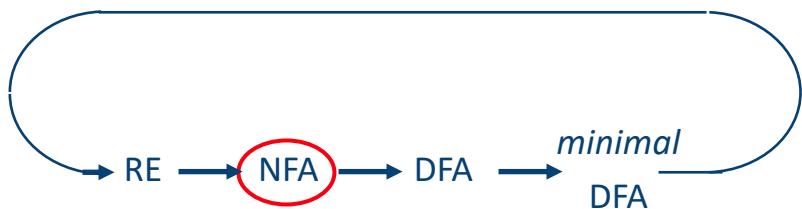


# Abbreviated Register Specification

Thompson's construction produces



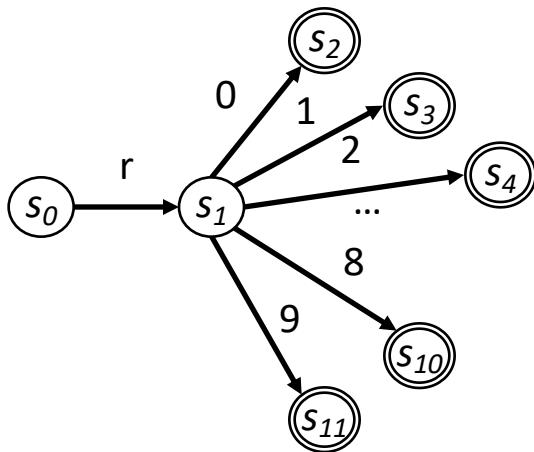
To make the example fit, we have  
eliminated some of the  $\epsilon$ -transitions, e.g.,  
between r and 0



# Abbreviated Register Specification



Applying the subset construction yields



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

The Cycle of Constructions

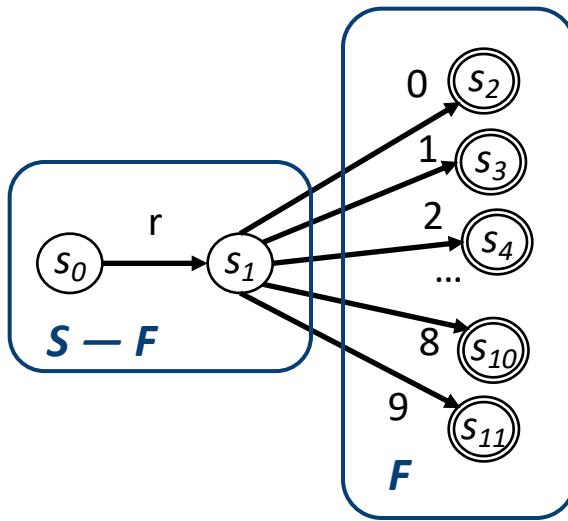


# Abbreviated Register Specification



## Hopcroft's algorithm

### Initial sets



$F$  does not split.

Since no transitions leave it, there are no states to split it.

### The Cycle of Constructions



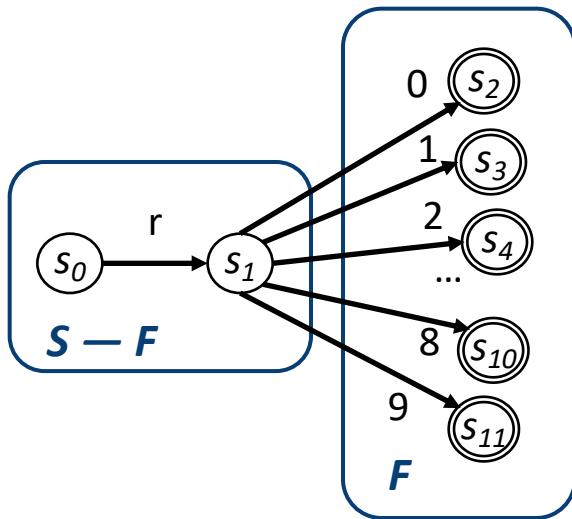
Technically, this edge shows up as 10 transitions, which are combined by construction of the character classifier ...

# Abbreviated Register Specification



## Hopcroft's algorithm

### Initial sets



$\{S - F\}$  does split

Any character in  $\Sigma$  will split it into  $\{s_0\}, \{s_1\}$

### The Cycle of Constructions



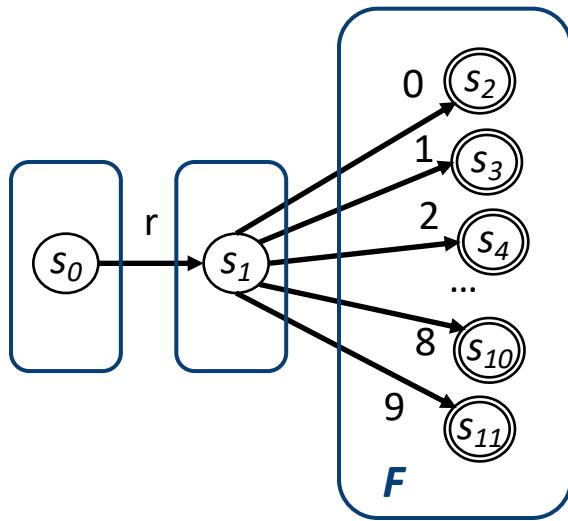
Technically, this edge shows up as 10 transitions, which are combined by construction of the character classifier ...

# Abbreviated Register Specification



## Hopcroft's algorithm

### Initial sets



$\{S - F\}$  does split

Any character in  $\Sigma$  will split it into  $\{s_0\}, \{s_1\}$

This partition is the final partition

### The Cycle of Constructions



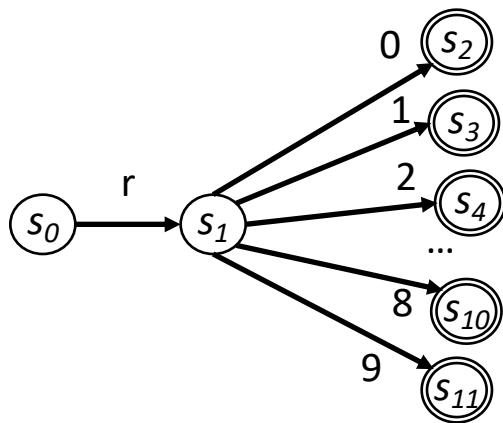
Technically, this edge shows up as 10 transitions, which are combined by construction of the character classifier ...

# Abbreviated Register Specification

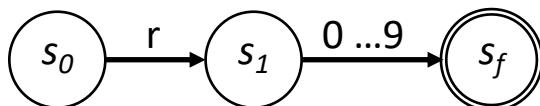


## Hopcroft's algorithm

Initial sets



Becomes, through minimization



### The Critical Takeaway Points:

- 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 “intuitive” RE

### The Cycle of Constructions





# 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

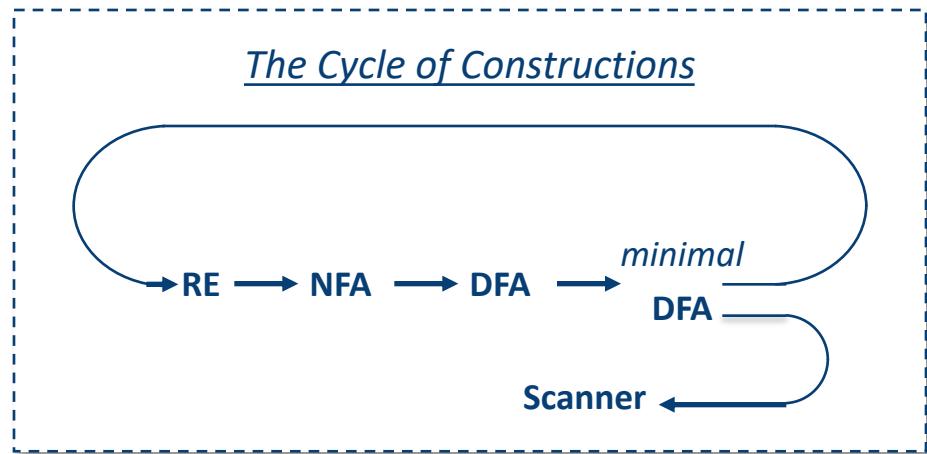
**Minimal DFA → Scanner**

- 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*

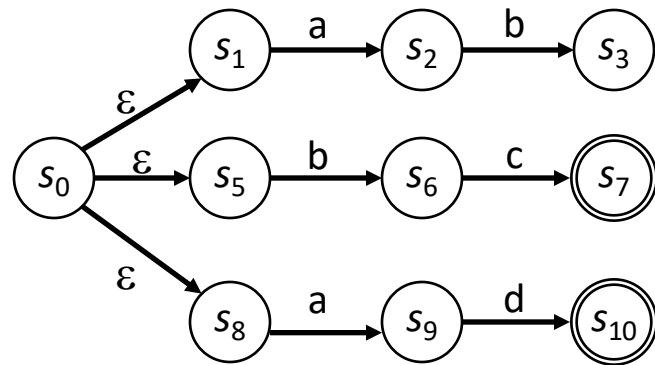


# Brzozowski's Algorithm for DFA Minimization



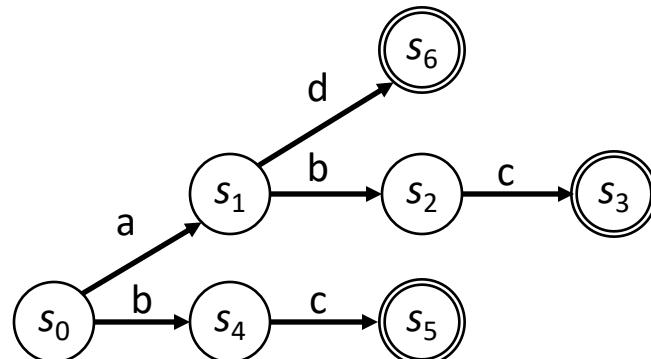
## The Intuition

- 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.*

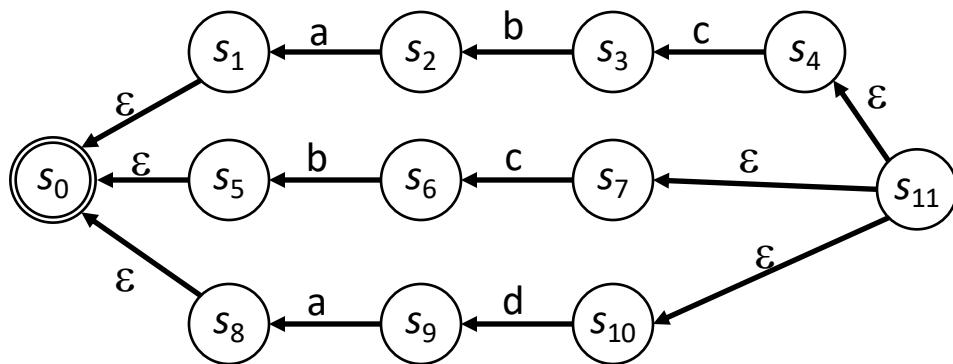
*Neither algorithm dominates the other.*



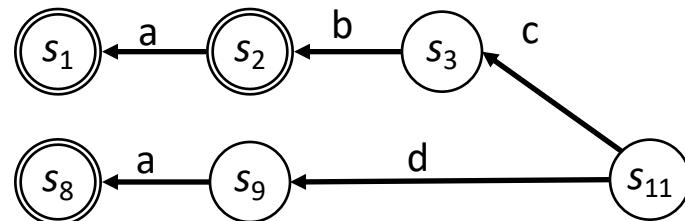
# Brzozowski's Algorithm

## Step 1

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



Reversed NFA



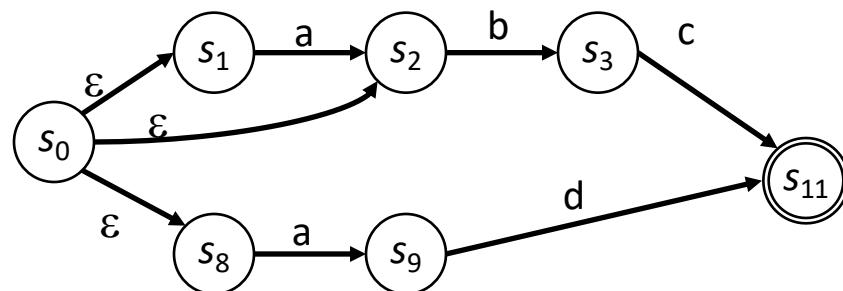
subset( $\text{reverse}(\text{NFA})$ )

# Brzozowski's Algorithm

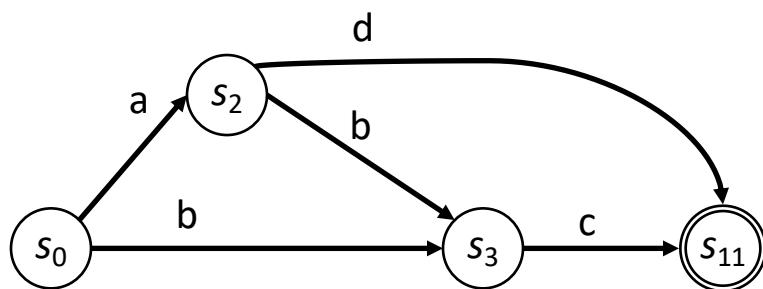


## Step 2

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



Reverse it, again



And subset it, again

The Cycle of Constructions

Minimal DFA



# Abbreviated Register Specification



Start with a regular expression

$r0 \mid r1 \mid r2 \mid r3 \mid r4 \mid r5 \mid r6 \mid r7 \mid r8 \mid r9$

Register names from zero to nine

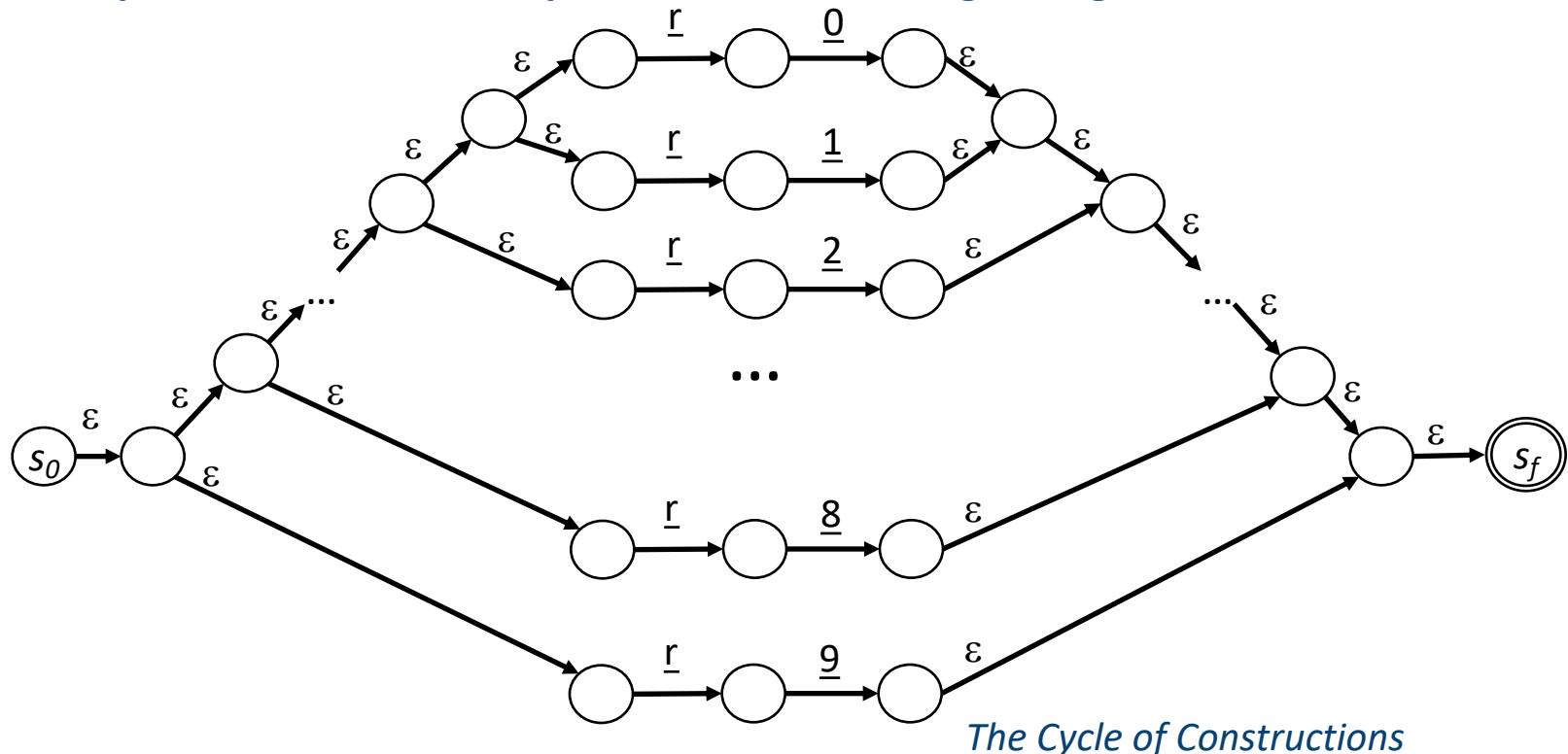
The Cycle of Constructions



# Abbreviated Register Specification



Thompson's construction produces something along these lines



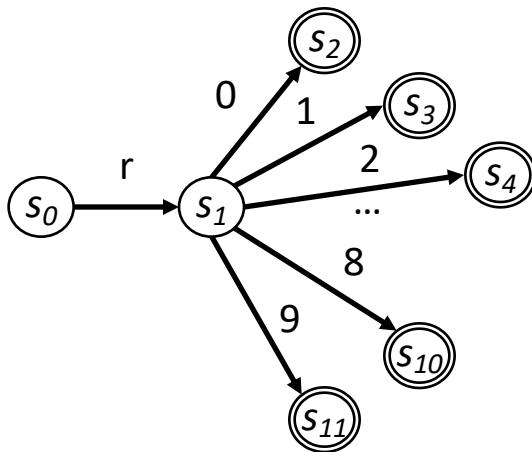
To make the example fit, we have  
eliminated some of the  $\epsilon$ -transitions, e.g.,  
between r and 0



# Abbreviated Register Specification



Applying the subset construction yields



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

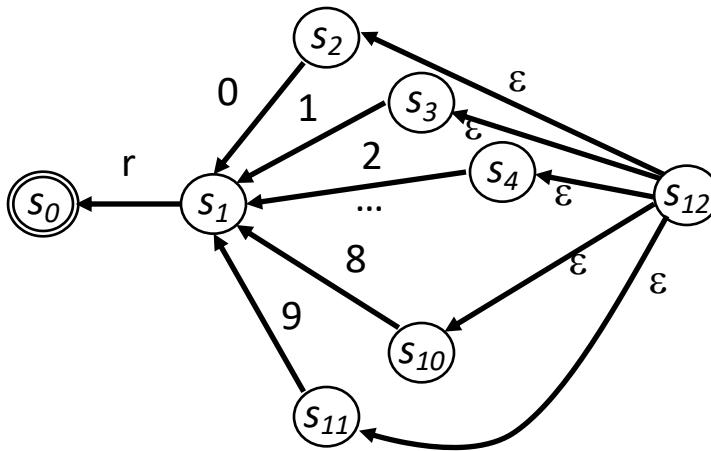
The Cycle of Constructions



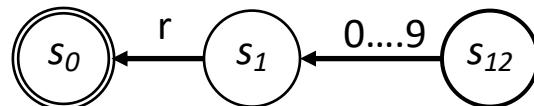
# Abbreviated Register Specification



Applying Brzozowski's algorithm, step 1



Reversed NFA



After Subset Construction

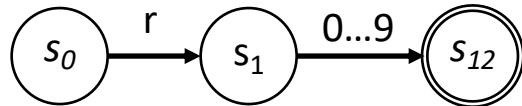
The Cycle of Constructions



# 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 “intuitive” **RE**

## Cycle of Constructions



# One Last Algorithm

## 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
             $R^0_{ii} = R^0_{ii} \cup \{\epsilon\};$ 

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

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

$R^k_{ij}$  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

