

Lexical Analysis: DFA Minimization

Comp 412

Copyright 2009, 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.

Automating Scanner Construction



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 (today)

- Hopcroft's algorithm

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

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

The Cycle of Constructions



DFA Minimization



The Big Picture

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

DFA Minimization



The Big Picture

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

Two states are equivalent if and only if:

- The set of paths leading to them are equivalent
- $\forall \alpha \in \Sigma$, transitions on α lead to equivalent states
- α -transitions to distinct sets \Rightarrow state must be split

(DFA)



DFA Minimization

The Big Picture

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

Two states are equivalent if and only if:

- The set of paths leading to them are equivalent
- $\forall \alpha \in \Sigma$, transitions on α lead to equivalent states (DFA)
- α -transitions to distinct sets \Rightarrow state must be split

A partition P of S

- A collection of sets P s.t. each $s \in S$ is in exactly one $p_i \in P$
- The algorithm iteratively constructs partitions of the DFA's states



DFA Minimization

Maximally sized sets \Rightarrow
minimal number of states

Details of the algorithm

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

Initial partition, P_0 , has two sets: $\{F\}$ & $\{S-F\}$ $D = (S, \Sigma, \delta, s_0, F)$
final states others

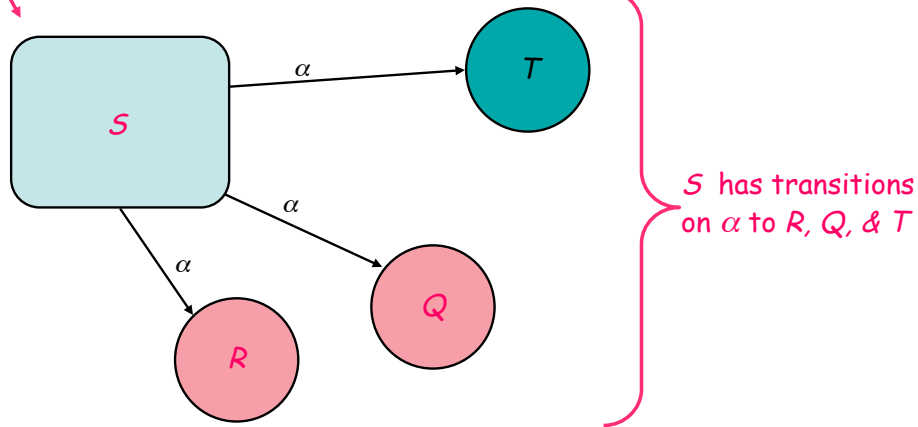
Splitting a set ("partitioning a set by \underline{a} ")

- Assume s_a & $s_b \in p_i$, and $\delta(s_a, \underline{a}) = s_x$, & $\delta(s_b, \underline{a}) = s_y$
- If s_x & s_y are not in the same set p_j , then p_i must be split
 - s_a has transition on \underline{a} , s_b does not $\Rightarrow \underline{a}$ splits p_i
- One state in the final DFA cannot have two transitions on \underline{a}



Key Idea: Splitting S around α

Original set S

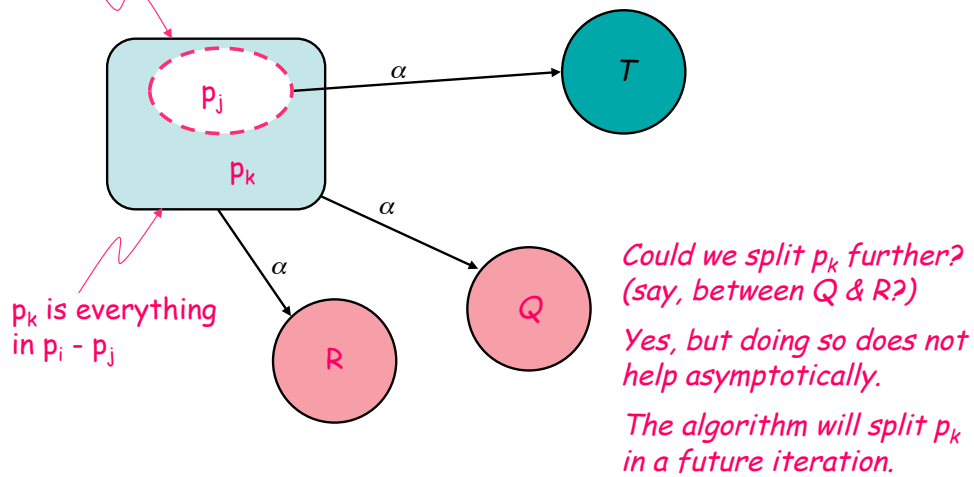


The algorithm partitions S around α



Key Idea: Splitting p_i around α

Original set p_i



This is a fixed-point algorithm!



DFA Minimization

The algorithm

```
T ← { F, {S-F} }
P ← { }
while ( P ≠ T )
  P ← T
  T ← { }
  for each set pi ∈ P
    T ← T ∪ Split(pi)

Split(S)
  for each c ∈ Σ
    if c splits S into s1 & s2
      then return {s1, s2}
  return S
```

Why does this work?

- Partition $P \in 2^S$
- Start off with 2 subsets of S : $\{F\}$ and $\{S-F\}$
- The *while* loop takes $P^i \rightarrow P^{i+1}$ by splitting 1 or more sets
- P^{i+1} is at least one step closer to the partition with $|S|$ sets
- Maximum of $|S|$ splits

Note that

- Partitions are never combined
- Initial partition ensures that final states remain final states

mild abuse of notation

DFA Minimization



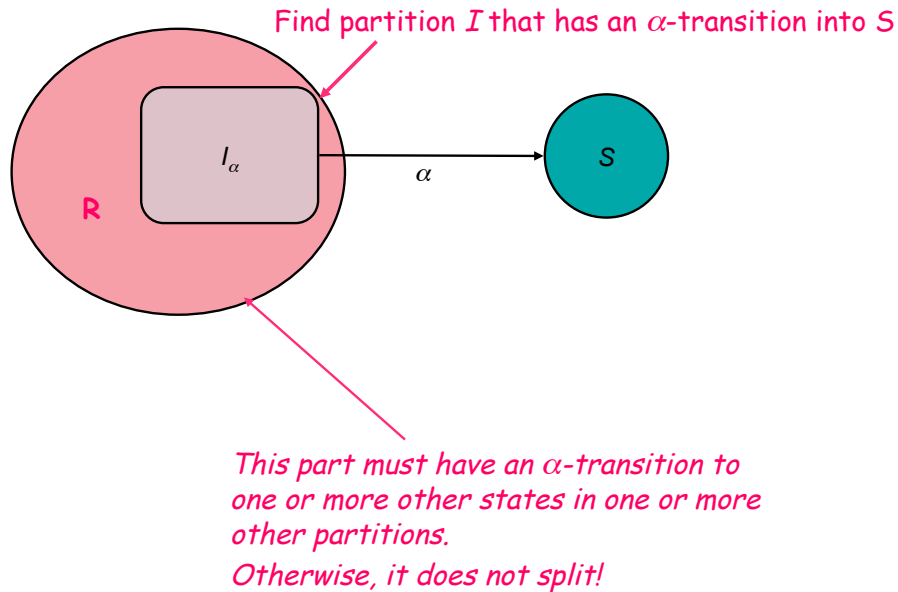
Refining the algorithm

- As written, it examines every $p_i \in P$ on each iteration
 - This strategy entails a lot of unnecessary work
 - Only need to examine p_i if some T , reachable from p_i , has split
- Reformulate the algorithm using a "worklist"
 - Start worklist with initial partition, F and $\{S-F\}$
 - When it splits p_i into p_1 and p_2 , place p_2 on worklist

This version looks at each $p_i \in P$ many fewer times

- Well-known, widely used algorithm due to John Hopcroft

Key Idea: Splitting S around α



Hopcroft's Algorithm



```

 $W \leftarrow \{F, S-F\}; P \leftarrow \{F, S-F\};$  //  $W$  is the worklist,  $P$  the current partition
while ( $W$  is not empty) do begin
  select and remove  $s$  from  $W$ ; //  $s$  is a set of states
  for each  $\alpha$  in  $\Sigma$  do begin
    let  $I_\alpha \leftarrow \delta_\alpha^{-1}(s)$ ; //  $I_\alpha$  is set of all states that can reach  $s$  on  $\alpha$ 
    for each  $p \in P$  such that  $p \cap I_\alpha$  is not empty
      and  $p$  is not contained in  $I_\alpha$  do begin
        partition  $p$  into  $p_1$  and  $p_2$  such that  $p_1 \leftarrow p \cap I_\alpha$ ;  $p_2 \leftarrow p - p_1$ ;
         $P \leftarrow (P - p) \cup p_1 \cup p_2$ ;
        if  $p \in W$ 
          then  $W \leftarrow (W - p) \cup p_1 \cup p_2$ ;
          else if  $|p_1| \leq |p_2|$ 
            then  $W \leftarrow W \cup p_1$ ;
            else  $W \leftarrow W \cup p_2$ ;
        end
      end
    end
  end

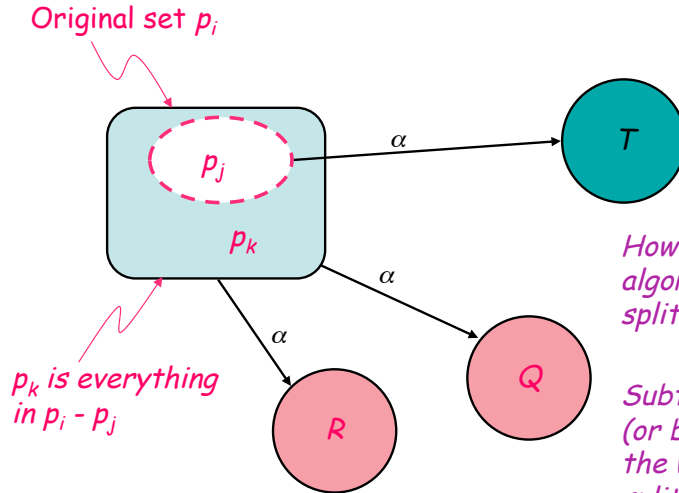
```

Critical difference between this formulation and the earlier one: this algorithm looks backward from a set; previously, it looked forward.

This distinction is critical to the worklist formulation. By projecting backward across the transitions, the algorithm can rely on the new partition to split its antecedents in the graph. This shows up in the example of a $(b|c)^*$ later in lecture.



Key Idea: Splitting p_i around α



Original set p_i

p_k is everything in $p_i - p_j$

How does the worklist algorithm ensure that it splits p_k around Q & R ?

Subtle point: either Q or R (or both) must already be on the worklist. (Q & R have split from $\{S-F\}$.)

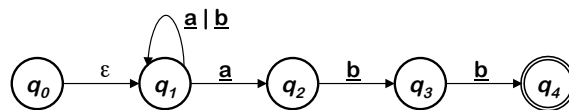
Thus, it can split p_i around one state (T) & add either p_j or p_k to the worklist.



A Detailed Example

Remember $(a | b)^* abb$?

(from last lecture)



Our first NFA

Applying the subset construction:

Iter.	State		ϵ -closure(move($s_i, *$))	
	DFA	NFA	\underline{a}	\underline{b}
0	s_0	q_0, q_1	q_1, q_2	q_1
1	s_1	q_1, q_2	q_1, q_2	q_1, q_3
	s_2	q_1	q_1, q_2	q_1
2	s_3	q_1, q_3	q_1, q_2	q_1, q_4
3	s_4	q_1, q_4	q_1, q_2	q_1

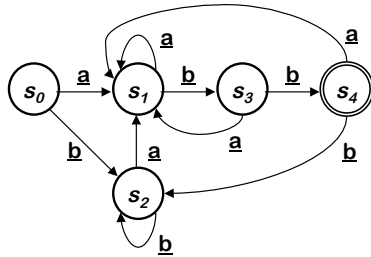
Iteration 3 adds nothing to S , so the algorithm halts

contains q_4 (final state)



A Detailed Example

The DFA for $(a | b)^* abb$



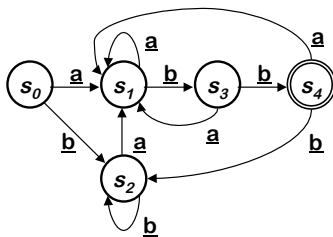
State	Character	
	a	b
s_0	s_1	s_2
s_1	s_1	s_3
s_2	s_1	s_2
s_3	s_1	s_4
s_4	s_1	s_2

- Not much expansion from NFA *(we feared exponential blowup)*
- Deterministic transitions
- Use same code skeleton as before



A Detailed Example (DFA Minimization)

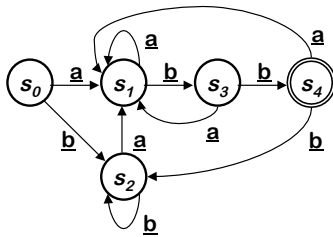
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$			



A Detailed Example (DFA Minimization)



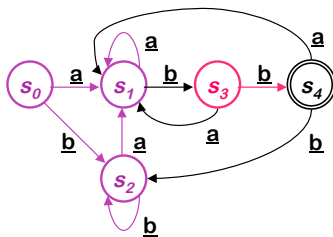
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	



A Detailed Example (DFA Minimization)



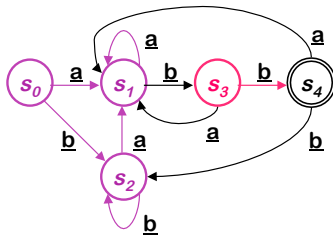
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\} \{s_0, s_1, s_2\}$



A Detailed Example (DFA Minimization)



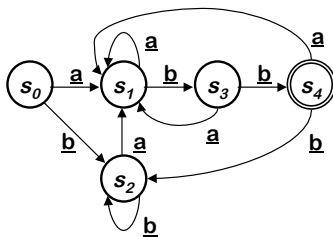
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\} \{s_0, s_1, s_2\}$
P_1	$\{s_4\} \{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\} \{s_0, s_1, s_2\}$			



A Detailed Example (DFA Minimization)



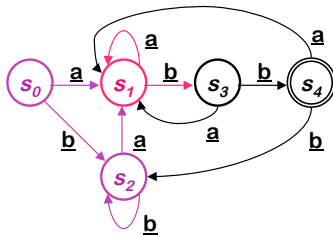
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\} \{s_0, s_1, s_2\}$
P_1	$\{s_4\} \{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\}$	none	



A Detailed Example (DFA Minimization)



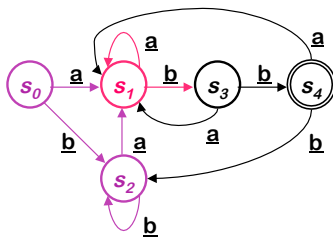
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\} \{s_0, s_1, s_2\}$
P_1	$\{s_4\} \{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\}$	none	$\{s_1\} \{s_0, s_2\}$



A Detailed Example (DFA Minimization)



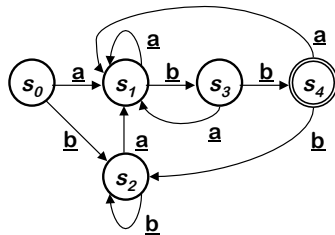
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\} \{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\} \{s_0, s_1, s_2\}$
P_1	$\{s_4\} \{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\} \{s_0, s_1, s_2\}$	$\{s_3\}$	none	$\{s_1\} \{s_0, s_2\}$
P_2	$\{s_4\} \{s_3\} \{s_1\} \{s_0, s_2\}$	$\{s_1\} \{s_0, s_2\}$			



A Detailed Example (DFA Minimization)



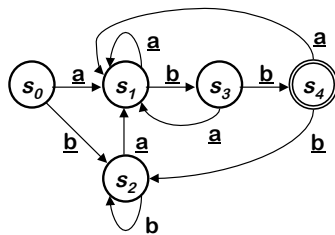
	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\}\{s_0,s_1,s_2,s_3\}$	$\{s_4\}\{s_0,s_1,s_2,s_3\}$	$\{s_4\}$	none	$\{s_3\}\{s_0,s_1,s_2\}$
P_1	$\{s_4\}\{s_3\}\{s_0,s_1,s_2\}$	$\{s_3\}\{s_0,s_1,s_2\}$	$\{s_3\}$	none	$\{s_1\}\{s_0,s_2\}$
P_2	$\{s_4\}\{s_3\}\{s_1\}\{s_0,s_2\}$	$\{s_1\}\{s_0,s_2\}$	$\{s_1\}$	none	none



A Detailed Example (DFA Minimization)



	Current Partition	Worklist	s	Split on <u>a</u>	Split on <u>b</u>
P_0	$\{s_4\}\{s_0,s_1,s_2,s_3\}$	$\{s_4\}\{s_0,s_1,s_2,s_3\}$	$\{s_4\}$	none	$\{s_3\}\{s_0,s_1,s_2\}$
P_1	$\{s_4\}\{s_3\}\{s_0,s_1,s_2\}$	$\{s_3\}\{s_0,s_1,s_2\}$	$\{s_3\}$	none	$\{s_1\}\{s_0,s_2\}$
P_2	$\{s_4\}\{s_3\}\{s_1\}\{s_0,s_2\}$	$\{s_1\}\{s_0,s_2\}$	$\{s_1\}$	none	none
P_2	$\{s_4\}\{s_3\}\{s_1\}\{s_0,s_2\}$	$\{s_1\}\{s_0,s_2\}$	$\{s_0,s_2\}$	none	none

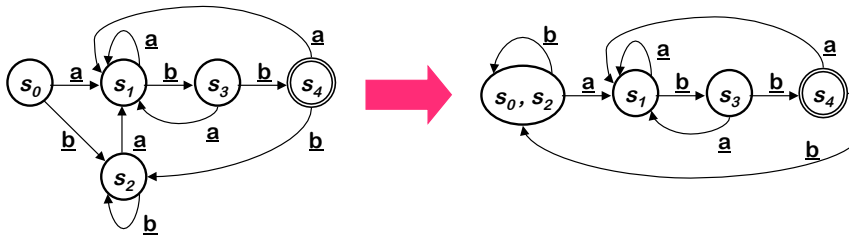


Empty worklist \Rightarrow done!

A Detailed Example (DFA Minimization)



	Current Partition	Worklist	s	Split on a	Split on b
P_0	$\{s_4\}\{s_0, s_1, s_2, s_3\}$	$\{s_4\}\{s_0, s_1, s_2, s_3\}$	$\{s_4\}$	none	$\{s_3\}\{s_0, s_1, s_2\}$
P_1	$\{s_4\}\{s_3\}\{s_0, s_1, s_2\}$	$\{s_3\}\{s_0, s_1, s_2\}$	$\{s_3\}$	none	$\{s_1\}\{s_0, s_2\}$
P_2	$\{s_4\}\{s_3\}\{s_1\}\{s_0, s_2\}$	$\{s_1\}\{s_0, s_2\}$	$\{s_1\}$	none	none
P_2	$\{s_4\}\{s_3\}\{s_1\}\{s_0, s_2\}$	$\{s_1\}\{s_0, s_2\}$	$\{s_0, s_2\}$	none	none



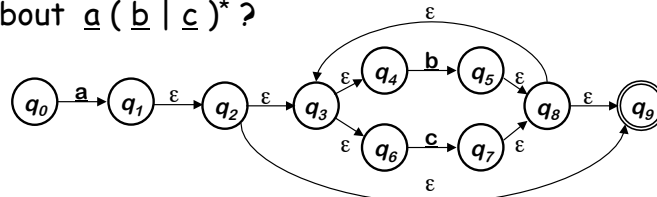
Comp 412, Fall 2009

20% reduction in number of states 24

DFA Minimization

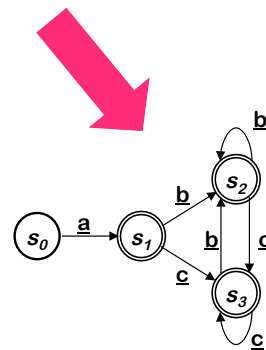


What about $\underline{a}(\underline{b} \mid \underline{c})^*$?



First, the subset construction:

States		ϵ -closure(Move($s, *$))		
DFA	NFA	a	b	c
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



From last lecture ...

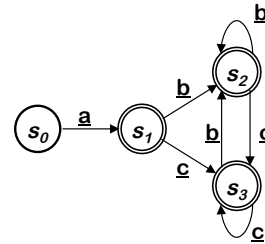
25

DFA Minimization



Then, apply the minimization algorithm

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



It splits no states after the initial partition

⇒ The minimal DFA has two states

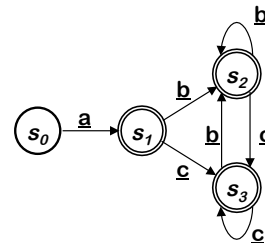
- One for $\{s_0\}$
- One for $\{s_1, s_2, s_3\}$

DFA Minimization

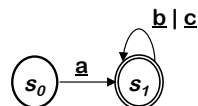


Then, apply the minimization algorithm

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



It produces this DFA



In lecture 5, we observed that a human would design a simpler automaton than Thompson's construction & the subset construction did.

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



Abbreviated Register Specification

Start with a regular expression

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

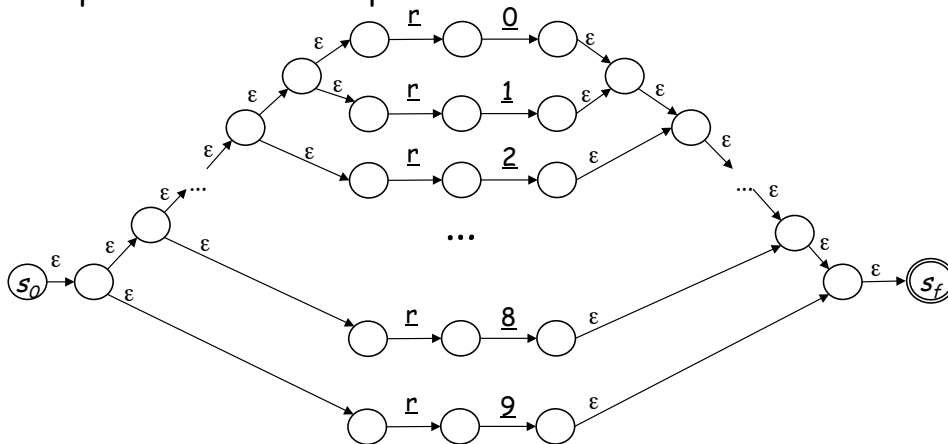
Register names from zero to nine

The Cycle of Constructions

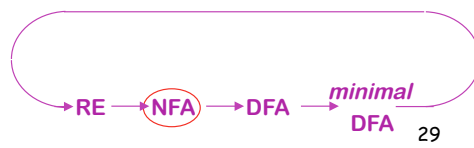


Abbreviated Register Specification

Thompson's construction produces



The Cycle of Constructions

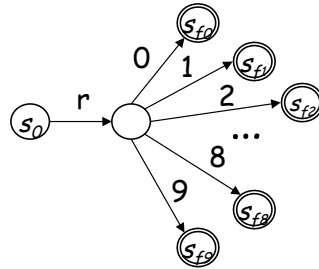


To make the example fit, we have eliminated some of the ϵ -transitions, e.g., between r and 0

Abbreviated Register Specification



The subset construction builds



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

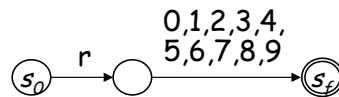
The Cycle of Constructions



Abbreviated Register Specification



The DFA minimization algorithm builds



This looks like what a skilled compiler writer would do!

The Cycle of Constructions

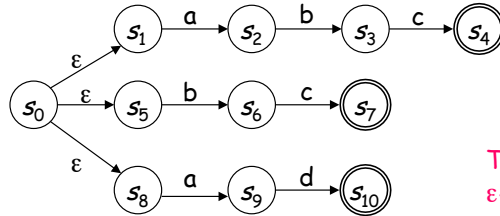




Alternative Approach to DFA Minimization

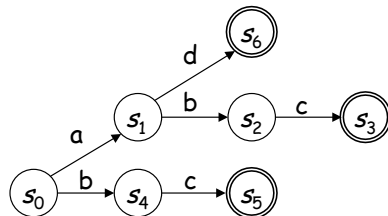
The Intuition

- 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.

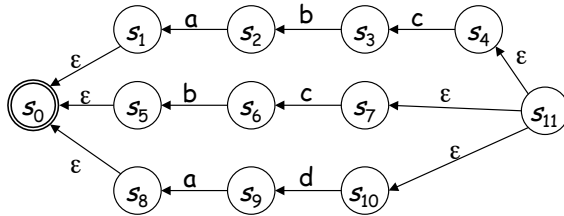
Neither algorithm dominates the other.



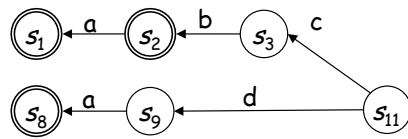
Alternative Approach to DFA Minimization

Step 1

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



Reversed NFA



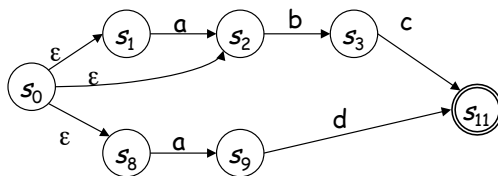
subset(reverse(NFA))



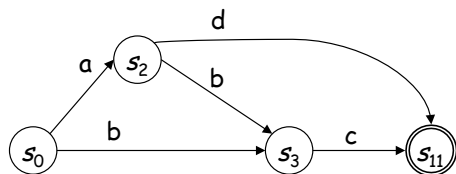
Alternative Approach to DFA Minimization

Step 2

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



Reverse it, again



Minimal DFA

And subset it, again

The Cycle of Constructions





RE Back to DFA

Kleene's Construction

```

for i ← 0 to |D| - 1;           // label each immediate path
  for j ← 0 to |D| - 1;
    R0ij ← { a | δ(di, a) = dj };
    if (i = j) then
      R0ii = R0ii | {ε};
  for k ← 0 to |D| - 1;       // label nontrivial paths
    for i ← 0 to |D| - 1;
      for j ← 0 to |D| - 1;
        Rkij ← Rk-1ik (Rk-1kk)* Rk-1kj | Rk-1ij
  L ← {}                       // union labels of paths from
  For each final state si     // s0 to a final state si
    L ← L | R|D|-10i

```

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

The Cycle of Constructions



Limits of Regular Languages

Not all languages are regular

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

You cannot construct DFA's to recognize these languages

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

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

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

- Strings with alternating 0's and 1's
($\epsilon \mid 1$)(01)^{*}($\epsilon \mid 0$)
- Strings with an even number of 0's and 1's

RE's can count bounded sets and bounded differences

Limits of Regular Languages



Advantages of Regular Expressions

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

Example — an expression grammar

$Term \rightarrow [a-zA-Z] ([a-zA-Z] | [0-9])^*$

$Op \rightarrow + | - | * | /$

$Expr \rightarrow (Term Op)^* Term$

Of course, this would generate a DFA ...

If REs are so useful ...

Why not use them for everything?