




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*) 

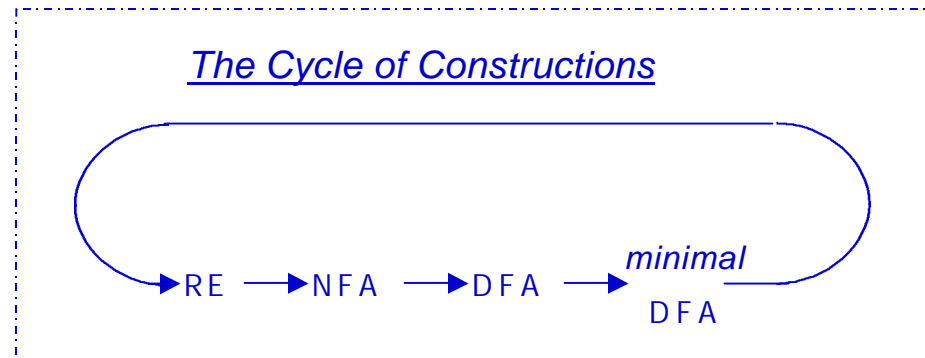
- Build the simulation

DFA \rightarrow Minimal DFA (*today*)

- Hopcroft's algorithm

DFA \rightarrow RE

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





DFA Minimization

The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one 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)
- transitions to distinct sets \Rightarrow states must be in distinct sets

A partition P of S

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



DFA Minimization

Details of the algorithm

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

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

Splitting a set

- Assume $q_a, q_b, \& q_c \in s$, and
- $\delta(q_a, \underline{a}) = q_x, \delta(q_b, \underline{a}) = q_y, \& \delta(q_c, \underline{a}) = q_z$
- If $q_x, q_y, \& q_z$ are not in the same set, then s must be split
- One state in the final DFA cannot have two transitions on \underline{a}



DFA Minimization

The algorithm

```

P ← { F, {Q-F} }
while ( P is still changing )
  T ← { }
  for each set s ∈ P
    for each a ∈ S
      partition s by a
      into s1, s2, ..., sk
      T ← T ∪ { s1, s2, ..., sk }
  if T ≠ P then
    P ← T

```

This is a fixed-point algorithm!

Why does this work?

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

Note that

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



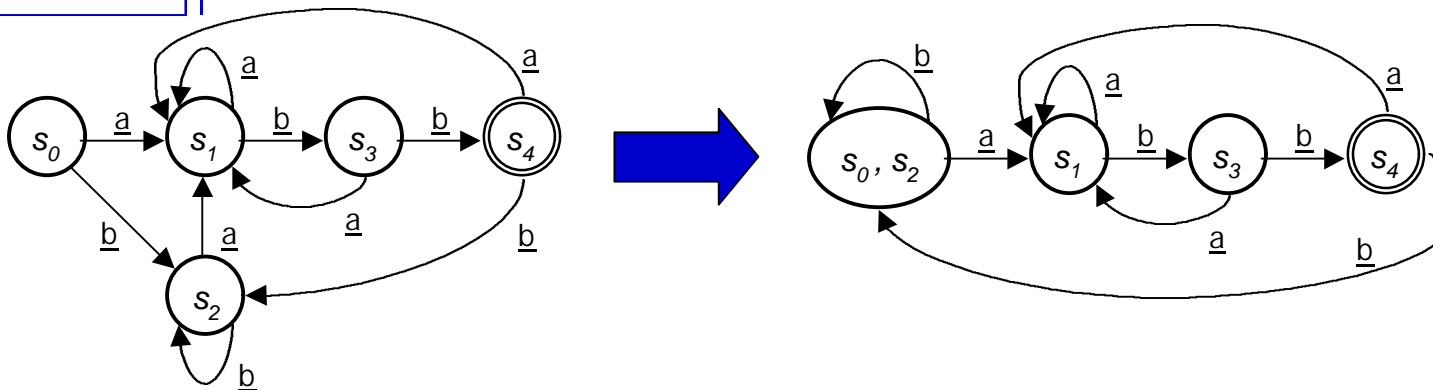
DFA Minimization

Enough theory, does this stuff work?

> Recall our example: $(\underline{a} \mid \underline{b})^* \underline{a}\underline{b}\underline{b}$

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

final state

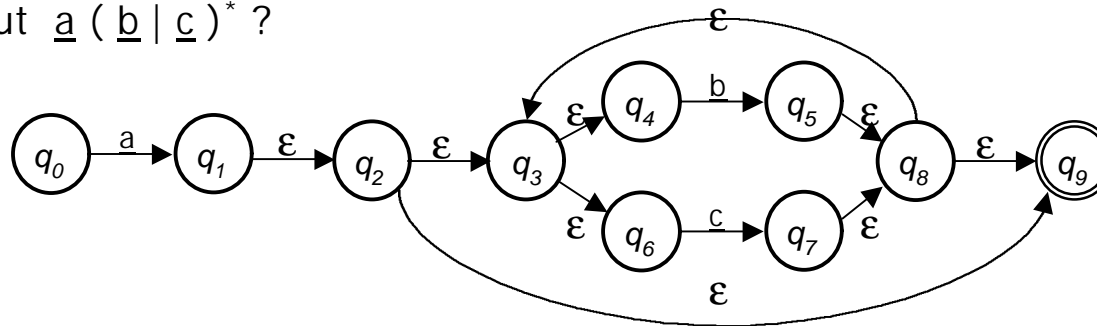


from Cooper & Torczon



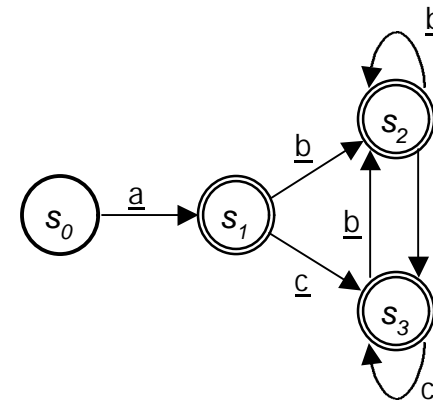
DFA Minimization

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



First, the subset construction:

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



Final states

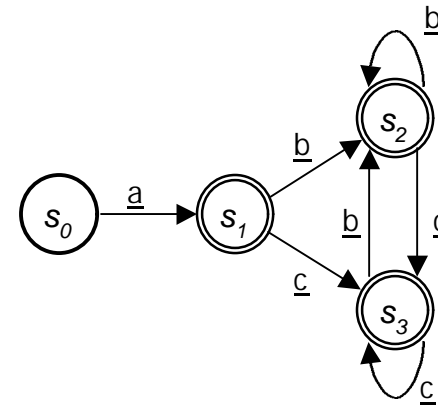
from Cooper & Torczon



DFA Minimization

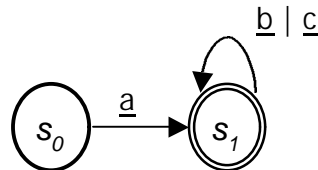
Then, apply the minimization algorithm

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



To produce the minimal DFA

final states



In lecture 6, I said that a human would design a simpler automaton than Thompson's construction did.
The algorithms produce that same DFA!



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

$Term \text{ @ } [a-zA-Z] ([a-zA-Z] | [\underline{0}-\underline{9}])^*$

$Op \text{ @ } + | - | * | /$

$Expr \text{ @ } (Term Op)^* Term$

Of course, this would generate a DFA ...

If REs are so useful ...

Why not use them for everything?



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)*
- $L = \{ w c w^r \mid w \in \Sigma^* \}$

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

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

- Alternating 0's and 1's
 $(\epsilon \mid 1)(0 \mid 1)(0 \mid \epsilon)$
- Sets of pairs of 0's and 1's
 $(01 \mid 10)(01 \mid 10)^*$

RE's can count bounded sets and bounded differences



What can be so hard?

Poor language design can complicate scanning

- Reserved words are important
if then then then = else; else else = then (PL/I)
- Significant blanks (Fortran & Algol68)
do 10 i = 1,25
do 10 i = 1.25
- String constants with special characters (C, others)
newline, tab, quote, comment delimiters, ...
- Finite closures
 - > Limited identifier length
 - > Adds states to count length

What can be so hard?

(Fortran 66/77)



```
      INTEGER FUNCTION A
      PARAMETER(A = 6, B = 2)
      IMPLICIT CHARACTER*(A-B) (A-B)
      INTEGER FORMAT(10), IF(10), DO9E1
100   FORMAT(4H) = (3)
200   FORMAT(4) = (3)
      DO9E1 = 1
      DO9E1 = 1, 2
   9   IF(X) = 1
      IF(X)H = 1
      IF(X)300, 200
300   CONTINUE
      END
C     THIS IS A "COMMENT CARD"
$   FILE(1)
      END
```

How does a compiler do this?

- First pass finds & inserts blanks
- Can add extra words or tags to create a scannable language
- Second pass is normal scanner

Example due to Dr. F.K. Zadeck