

CS415 Compilers

Lexical Analysis Part 4

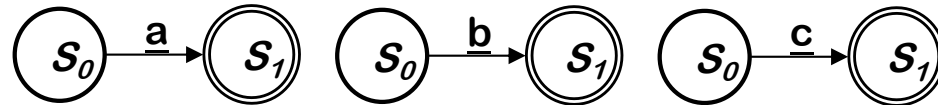
These slides are based on slides copyrighted by
Keith Cooper, Ken Kennedy & Linda Torczon at Rice
University

- Third homework has been posted. Due Monday, March 7
- First project (local instruction scheduler) NEW deadlines:
code: March 9
report: March 11

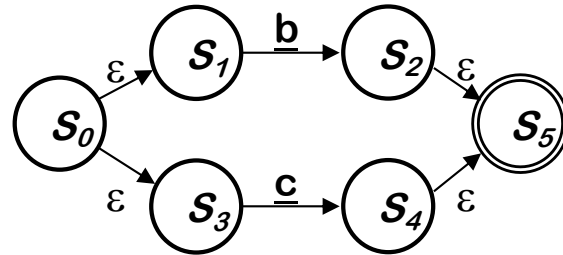
1. The order in which programs perform I/O has to be preserved. This requires a dependence notion on outputAI instructions.
2. Benchmark codes have storeAI and loadAI memory accesses only, with r0 as the base register. You must use the offset to determine whether a dependence exists or not.
3. We will have some private tests that have store and load memory accesses. This is for [extra credit](#).
4. Project reports should be around 6 pages long, with a max of 8 pages. We will not read your report beyond 8 pages. The report should include a short description of what you did, the outcome of your experiments, and how you interpret these outcomes. Use graphs/figures to show your results.

Let's try $a(b|c)^*$

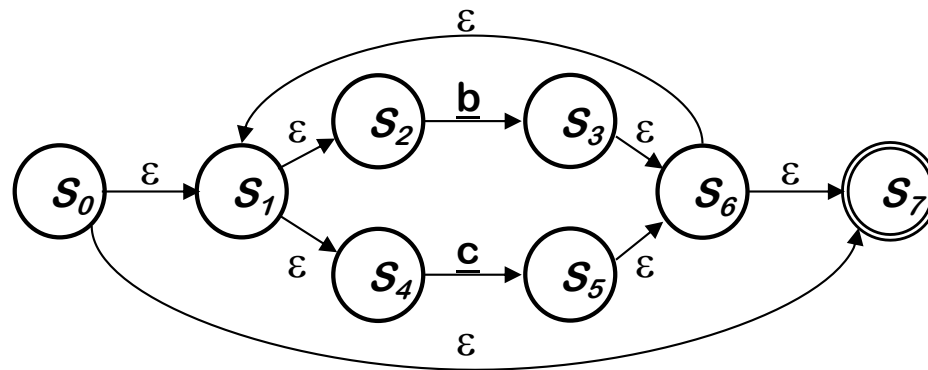
1. a , b , & c



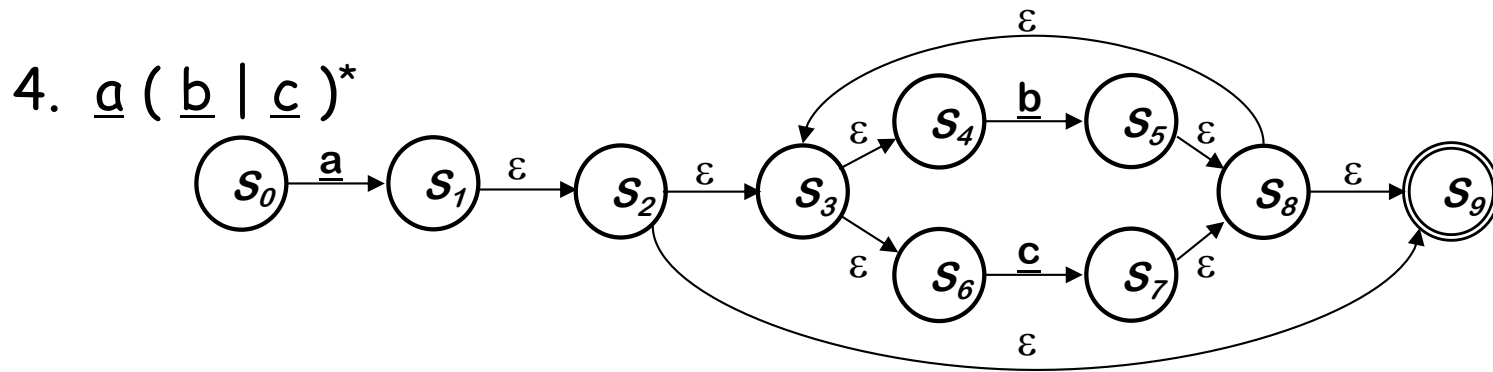
2. $b|c$



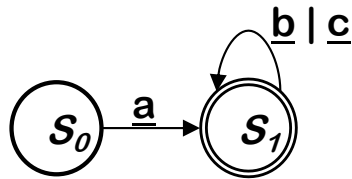
3. $(b|c)^*$



RUTGERS Example of Thompson's Construction (*con't*)



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



But, we can automate production of the more complex one ...

Need to build a simulation of the NFA

Two key functions

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

The algorithm (sketch):

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

Sounds more complex than it is...

The algorithm:

```

 $s_0 \leftarrow \varepsilon\text{-closure}(q_0)$ 
add  $s_0$  to  $S$ 
while (  $S$  is still changing )
  for each  $s_i \in S$ 
    for each  $a \in \Sigma$ 
       $s_j \leftarrow \varepsilon\text{-closure}(\text{move}(s_i, a))$ 
      if (  $s_j \notin S$  ) then
        add  $s_j$  to  $S$  as  $s_j$ 
         $T[s_i, a] \leftarrow s_j$ 
      else
         $T[s_i, a] \leftarrow s_j$ 

```

Let's think about why this works

The algorithm halts:

1. S contains no duplicates (test before adding)
2. 2^Q 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 symbol in each s_i .

It builds every possible NFA configuration.

$\Rightarrow S$ and T form the DFA

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

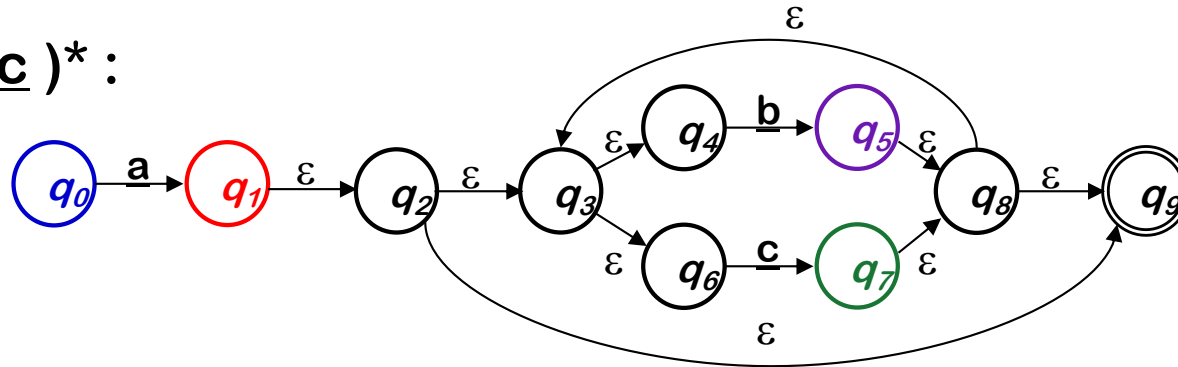
Other fixed-point computations

- Canonical construction of sets of LR(1) items
 - \rightarrow Quite similar to the subset construction
- Classic data-flow analysis
 - \rightarrow Solving sets of simultaneous set equations
- DFA minimization algorithm (coming up!)

We will see many more fixed-point computations

RUTGERS NFA \rightarrow DFA with Subset Construction

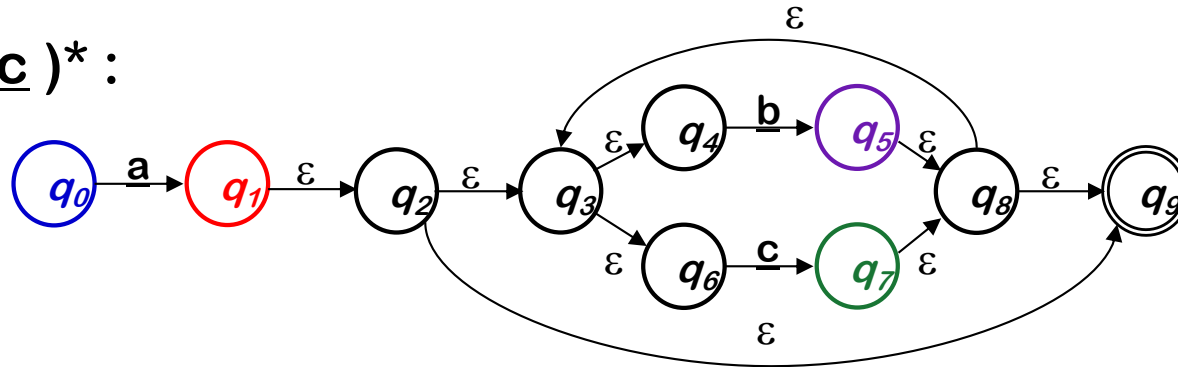
a (b | c)^{*} :



Applying the subset construction:

RUTGERS NFA \rightarrow DFA with Subset Construction

a (b | c)^{*} :

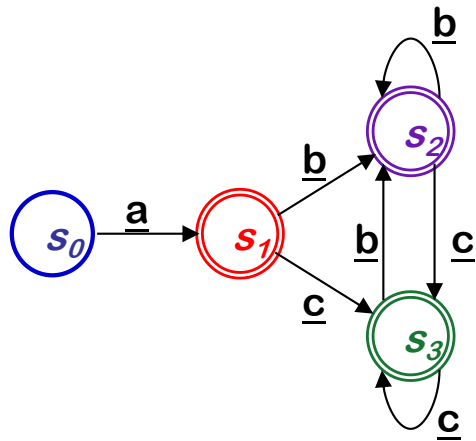


Applying the subset construction:

	Sets of NFA states	ϵ -closure(move(s,*))		
		<u>a</u>	<u>b</u>	<u>c</u>
<i>s</i> ₀	{ <i>q</i> ₀ }	{ <i>q</i> ₁ , <i>q</i> ₂ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ , <i>q</i> ₉ }	<i>none</i>	<i>none</i>
<i>s</i> ₁	{ <i>q</i> ₁ , <i>q</i> ₂ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ , <i>q</i> ₉ }	<i>none</i>	{ <i>q</i> ₅ , <i>q</i> ₈ , <i>q</i> ₉ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ }	{ <i>q</i> ₇ , <i>q</i> ₈ , <i>q</i> ₉ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ }
<i>s</i> ₂	{ <i>q</i> ₅ , <i>q</i> ₈ , <i>q</i> ₉ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ }	<i>none</i>	<i>s</i> ₂	<i>s</i> ₃
<i>s</i> ₃	{ <i>q</i> ₇ , <i>q</i> ₈ , <i>q</i> ₉ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆ }	<i>none</i>	<i>s</i> ₂	<i>s</i> ₃

Final states

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



δ	<u>a</u>	<u>b</u>	<u>c</u>
s_0	s_1	-	-
s_1	-	s_2	s_3
s_2	-	s_2	s_3
s_3	-	s_2	s_3

- Ends up smaller than the NFA
- All transitions are deterministic

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

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

NFA \rightarrow DFA (*subset construction*)

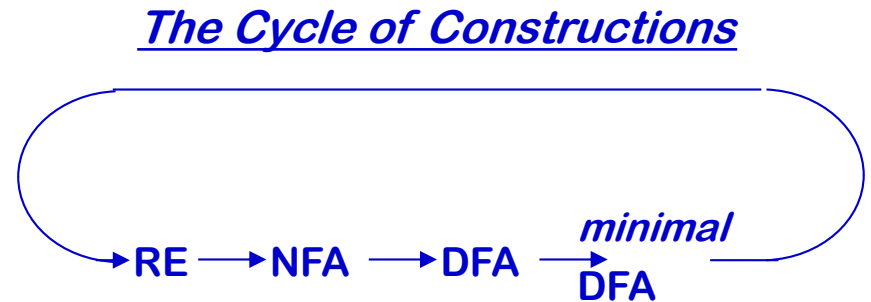
- Build the simulation

DFA \rightarrow Minimal DFA

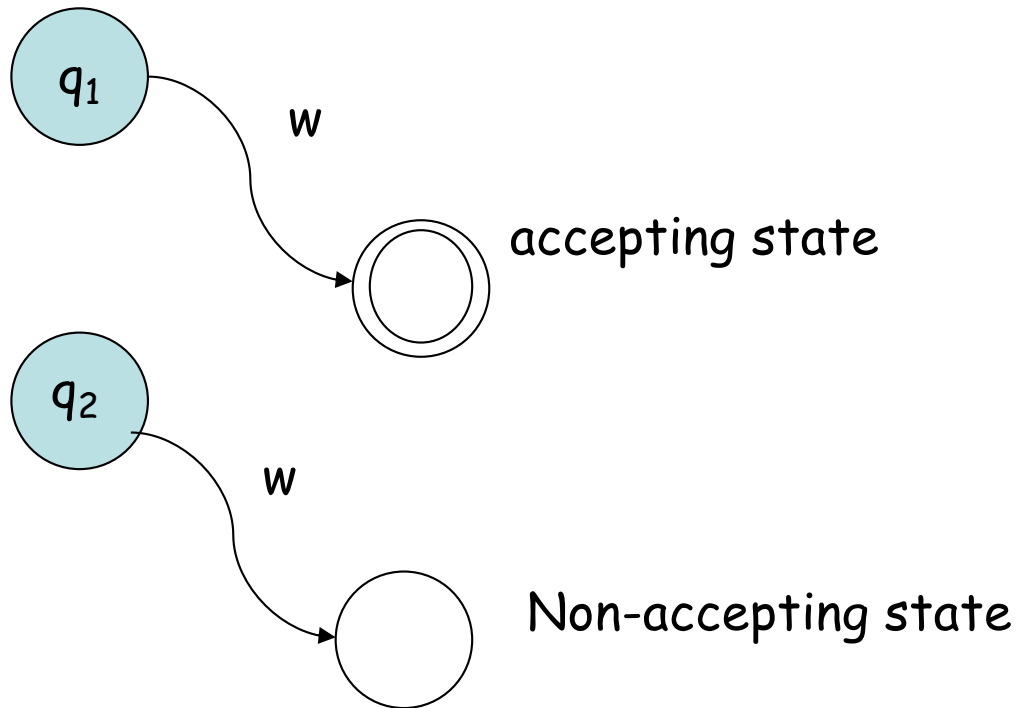
- Hopcroft's algorithm

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

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



How do we know whether two states
encode the same information?



Intuition: Two states are
equivalent if for all
sequences of input
symbols "w" they both
lead to an accepting state,
or both end up in a non-
accepting state.

q_1 and q_2 are not equivalent.
"w" is a witness that they
are not equivalent.

The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

Two states are equivalent if and only if:

- $\forall a \in \Sigma$, transitions on a lead to equivalent states (DFA)
- if a -transitions to different sets \Rightarrow two states must be in different sets, i.e., cannot be equivalent

The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

Two states are equivalent if and only if:

- $\forall a \in \Sigma$, transitions on a lead to equivalent states (DFA)
- if a -transitions to different sets \Rightarrow two states must be in different sets, i.e., cannot be equivalent

A partition P of S

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

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: $\{F\}$ & $\{Q-F\}$ ($D = (Q, \Sigma, \delta, q_0, F)$)

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

- Assume $q_a, \& \ q_b \in s$, and $\delta(q_a, \underline{a}) = q_x, \& \ \delta(q_b, \underline{a}) = q_y$
- If $q_x \& \ q_y$ are not in the same set, i.e., are considered equivalent, then s must be split
 - q_a has transition on \underline{a} , q_b does not $\Rightarrow \underline{a}$ splits s

The algorithm

```
 $P \leftarrow \{ F, \{Q-F\} \}$   
while ( $P$  is still changing)  
   $T \leftarrow \{ \}$   
  for each set  $S \in P$   
     $T \leftarrow T \cup \text{split}(S)$   
   $P \leftarrow T$   
  
split( $S$ ):  
  for each  $a \in \Sigma$   
    if  $a$  splits  $S$  into  
       $S_1, S_2, \dots$  then  
        return  $\{S_1, S_2, \dots\}$   
  else return  $S$ 
```

This is a fixed-point algorithm!

Why does this work?

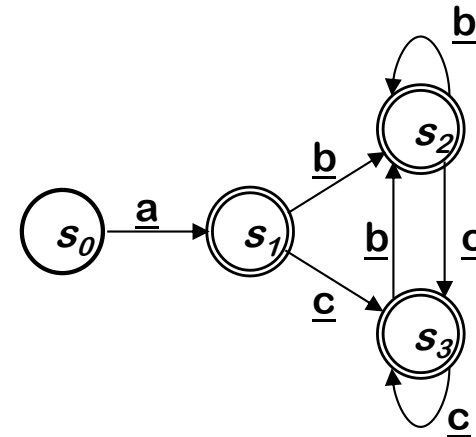
- Start off with 2 subsets of Q $\{F\}$ and $\{Q-F\}$
- *While* loop takes $P_i \rightarrow P_{i+1}$ by splitting 1 or more sets
- P_{i+1} is at least one step closer to the partition with $|Q|$ sets
- Maximum of $|Q|$ splits

Note that

- Partitions are never combined

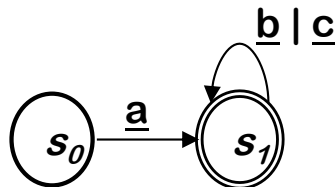
Then, apply the minimization algorithm

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



final states

To produce the minimal DFA



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!

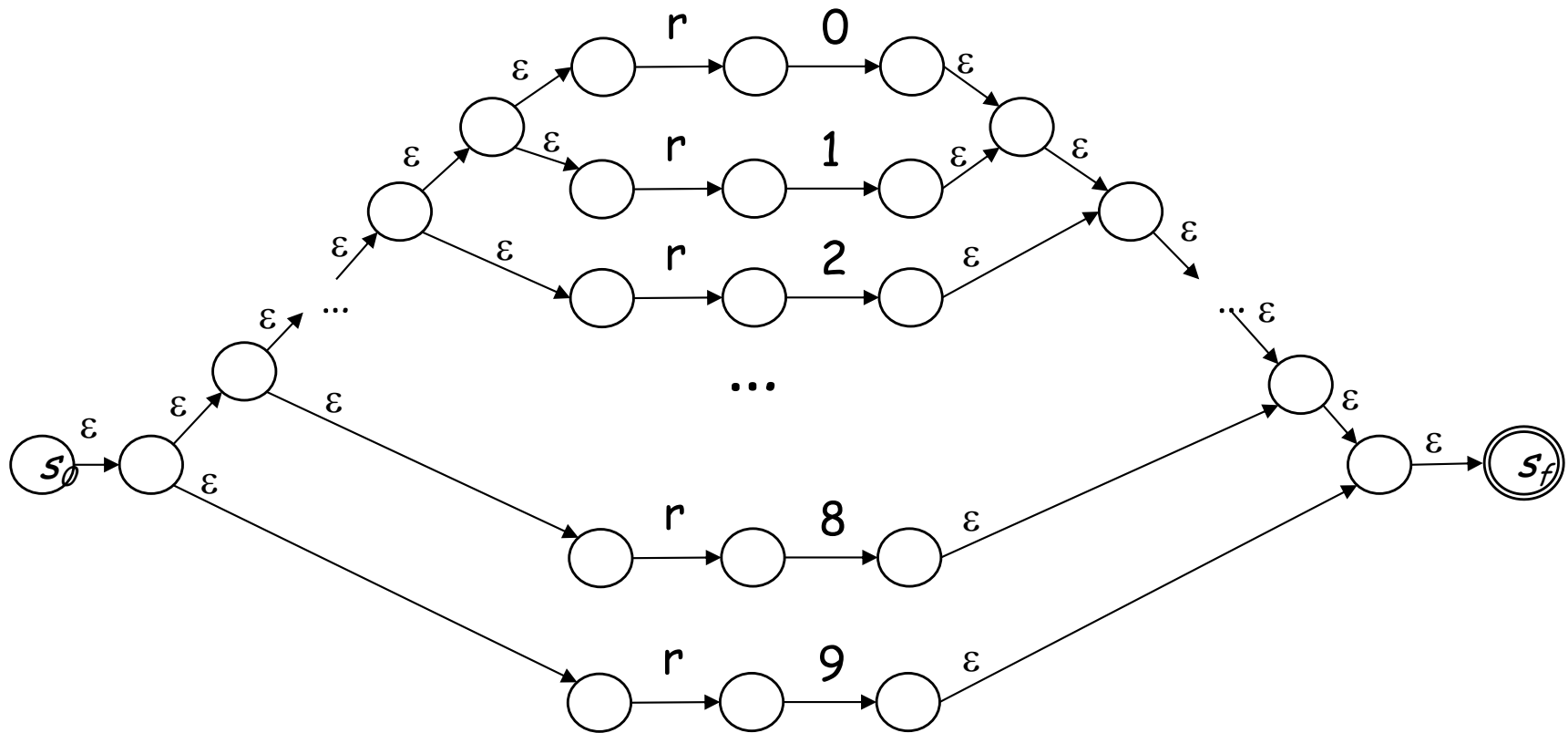
Start with a regular expression

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

The Cycle of Constructions



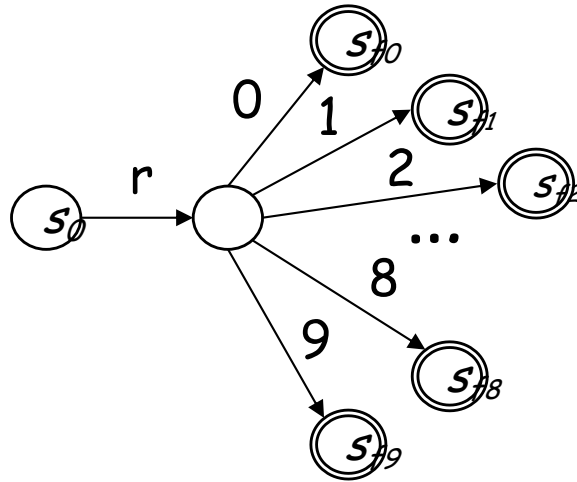
Thompson's construction produces



The Cycle of Constructions

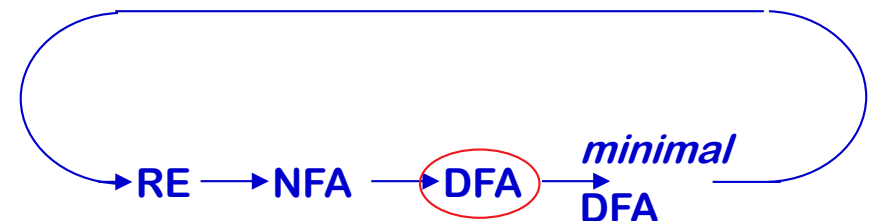


The subset construction builds

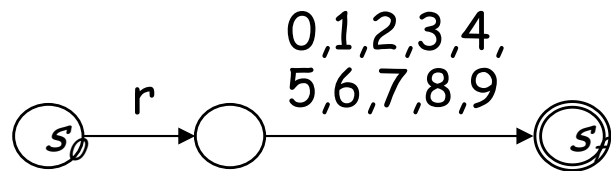


This is a DFA, but it has a lot of states ...

The Cycle of Constructions



The DFA minimization algorithm builds



This looks like what a skilled compiler writer would do!

The Cycle of Constructions



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

$$Term \rightarrow [a-zA-Z]([a-zA-z] | [\underline{0}-\underline{9}])^*$$
$$Op \rightarrow + | - | * | /$$
$$Expr \rightarrow (Term Op)^* Term$$

Of course, this would generate a DFA ...

If REs are so useful ...

Why not use them for everything?

Not all languages are regular

$$\text{RL's} \subset \text{CFL's} \subset \text{CSL's}$$

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

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

- Strings with alternating 0's and 1's
 $(\varepsilon \mid 1)(01)^*(\varepsilon \mid 0)$
- Strings with an even number of 0's and 1's
- Strings of bit patterns that represent binary numbers which are divisible by 5

Poor language design can complicate scanning

- Reserved words are important
if then then then = else; else else = then (PL/I)
- Insignificant blanks (Fortran & Algol68)
do 10 i = 1,25
do 10 i = 1.25
- String constants with special characters (C, C++, Java, ...)
newline, tab, quote, comment delimiters, ...
- Limited identifier "length" (Fortran 66 & PL/I)

Parsing (Syntax Analysis)

EAC Chapters 3.1 - 3.2