

# CSE 311 Foundations of Computing I

Lecture 24  
FSM Limits, Pattern Matching  
Autumn 2011

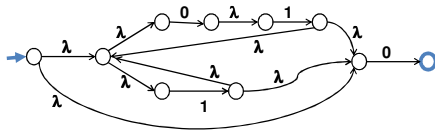
## Announcements

- Reading assignments
  - 7<sup>th</sup> Edition, Section 13.4
  - 6<sup>th</sup> Edition, Section 12.4
  - 5<sup>th</sup> Edition, Section 11.4

## Last lecture highlights

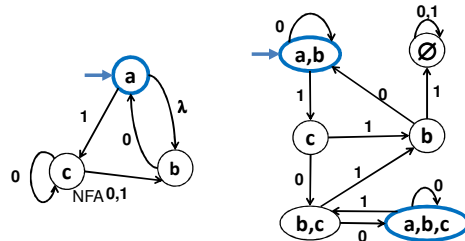
- NFAs from Regular Expressions

$(01 \cup 1)^*0$



## Last lecture highlights

- “Subset construction”: NFA to DFA



## What can Finite State Machines do?

- We’ve seen how we can get DFAs to recognize all regular languages
- What about some other languages we can generate with CFGs?
  - $\{0^n1^n : n \geq 0\}$ ?
  - Binary Palindromes?
  - Strings of Balanced Parentheses?

## $A = \{0^n1^n : n \geq 0\}$ cannot be recognized by any DFA

Consider the infinite set of strings

$S = \{\lambda, 0, 00, 000, 0000, \dots\}$

Claim: No two strings in  $S$  can end at the same state of any DFA for  $A$ , so no such DFA can exist

Proof: Suppose  $n \neq m$  and  $0^n$  and  $0^m$  end at the same state  $p$ .

Since  $0^n1^n$  is in  $A$ , following  $1^n$  after state  $p$  must lead to a final state.

But then the DFA would accept  $0^m1^n$  which is a contradiction

## The set B of binary palindromes cannot be recognized by any DFA

Consider the infinite set of strings

$S = \{\lambda, 0, 00, 000, 0000, \dots\}$

Claim: No two strings in S can end at the same state of any DFA for B, so no such DFA can exist

Proof: Suppose  $n \neq m$  and  $0^n$  and  $0^m$  end at the same state p.

Since  $0^n 10^n$  is in B, following  $10^n$  after state p must lead to a final state.

But then the DFA would accept  $0^m 10^n$  which is a contradiction

## The set P of strings of balanced parentheses cannot be recognized by any DFA

## Pattern Matching

- Given
  - a string, **s**, of **n** characters
  - a pattern, **p**, of **m** characters
  - usually  $m \ll n$
- Find
  - all occurrences of the pattern **p** in the string **s**
- Obvious algorithm:
  - try to see if **p** matches at each of the positions in **s**
    - stop at a failed match and try the next position

String **s** = x y x x y x y x y x y x y x y x y x y x x  
Pattern **p** = x y x y y x y x y x x

String **s** = x y x x y x y x y x y x y x y x y x y x x  
          x y x y y x y x y x x

String **s** = x y x x y x y x y x y x y x y x y x y x x  
          x y x y  
          x y x y y x y x y x x

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y x y y x y x y x x

13

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y  
x y x y y x y x y x x

14

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y  
x y x y y  
x y x y y x y x y x x

15

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y  
x y x y y  
x  
x y x y y x y x y x x

16

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y  
x y x y y  
x  
x y x y y x y x y x x  
x y x y y x y x y x x

17

String **s** = x y x x y x y x y y x y x y x y y x y x y x x  
x y x y  
x  
x y  
x y x y y  
x  
x y x y y x y x y x x  
x  
x y x y y x y x y x x

18

String  $s = x y x x y x y x y y x y x y x x$

```

x y x y
x
x y
x y x y y
x
x y x y y x y x y x x
x
x y x
x y x y y x y x y x x

```

19

String  $s = x y x x y x y x y y x y x y x y x x$

```

x y x y
x
x y
x y x y y
x
x y x y y x y x y x x
x
x y x
x y x y y x y x y x x

```

20

String  $s = x y x x y x y x y y x y x y x y x y x x$

```

x y x y
x
x y
x y x y y
x
x y x y y x y x y x x
x
x y x
x
x
x y x y y x y x y x x

```

21

String  $s = x y x x y x y x y y x y x y x y x y x x$

```

x y x y
x
x y
x y x y y
x
x y x y y x y x y x x
x
x y x
x
x
x y x y y
x y x y y x y x y x x

```

22

String  $s = x y x x y x y x y y x y x y x y x y x x$

```

x y x y
x
x y
x y x y y
x
x y x y y x y x y x x
x
x y x
x
x
x y x y y
x
x y x y y x y x y x x

```

Worst-case time  
 $O(mn)$

23

String  $s = x y x x y x y x y y x y x y x y x y x x$

```

x y x y
x y x y
x y x y y
x
x y x
x
x
x y x y y
x y x y y x y x y x x

```

Lots of wasted work

24

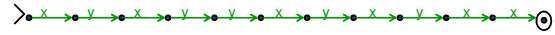
### Better Pattern Matching via Finite Automata

- Build a DFA for the pattern (preprocessing) of size  $O(m)$ 
  - Keep track of the 'longest match currently active'
  - The DFA will have only  $m+1$  states
- Run the DFA on the string  $n$  steps
- Obvious construction method for DFA will be  $O(m^2)$  but can be done in  $O(m)$  time.
- Total  $O(m+n)$  time

25

### Building a DFA for the pattern

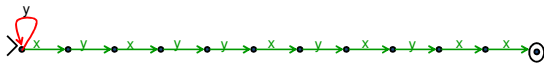
Pattern  $p = x y x y y x y x y x x$



26

### Preprocessing the pattern

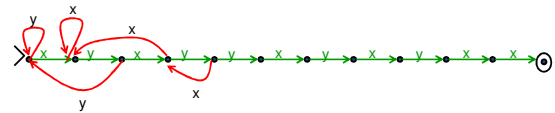
Pattern  $p = x y x y y x y x y x x$



27

### Preprocessing the pattern

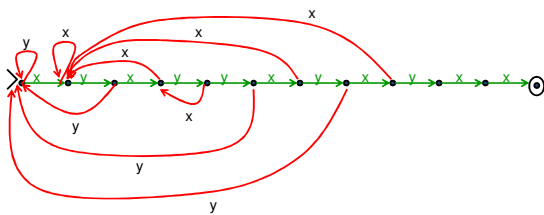
Pattern  $p = x y x y y x y x y x x$



28

### Preprocessing the pattern

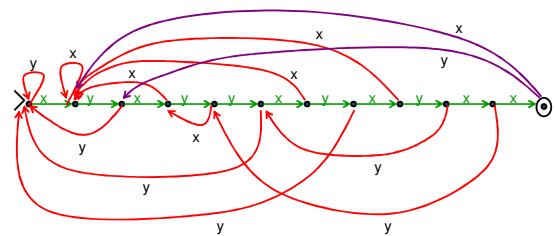
Pattern  $p = x y x y y x y x y x x$



29

### Preprocessing the pattern

Pattern  $p = x y x y y x y x y x x$



30

## Generalizing

- Can search for arbitrary combinations of patterns
  - Not just a single pattern
  - Build NFA for pattern then convert to DFA 'on the fly'.
    - Compare DFA constructed above with subset construction for the obvious NFA.

31

## A Quick Note...

- On how to convert NFAs and DFAs to equivalent regular expressions...
- We've already seen
  - DFAs and NFAs recognize the same languages
  - NFAs (and therefore DFAs) recognize any language given by a regular expression
- This completes the equivalence of DFAs and regular expressions

Autumn 2011

CSE 311

32

## Generalized NFAs

- Like NFAs but allow
  - Parallel edges
  - Regular Expressions as edge labels
    - NFAs already have edges labeled  $\lambda$  or  $\alpha$
- An edge labeled by **A** can be followed by reading a string of input chars that is in the language represented by **A**
- A string  $x$  is accepted iff there is a path from start to final state labeled by a regular expression whose language contains  $x$

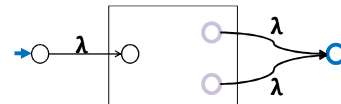
Autumn 2011

CSE 311

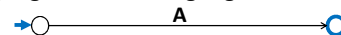
33

## Starting from NFA

- Add new start state and final state



- Then eliminate original states one by one, keeping the same language, until it looks like:



- Final regular expression will be **A**

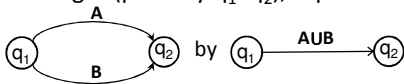
Autumn 2011

CSE 311

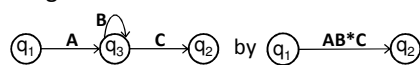
34

## Only two simplification rules:

- **Rule 1:** For any two states  $q_1$  and  $q_2$  with parallel edges (possibly  $q_1=q_2$ ), replace



- **Rule 2:** Eliminate non-start/final state  $q_3$  by replacing all



for every pair of states  $q_1, q_2$  (even if  $q_1=q_2$ )

Autumn 2011

CSE 311

35