

CSE 322: Introduction to Formal Models in Computer Science

Pattern Matching

Paul Beame

1

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 , stopping at a failed match

2

String $s = xyxyxyxyxyxyxyxyxyxyxyxyxy$
Pattern $p = xyxyxyxyxy$

3

String $s = xyxxyxyxyxyxyxyxyxyxyxyxyxy$
 $xyxyxyxyxy$

4

String $s = xyxyxyxyxyxyxyxyxyxyxyxyxy$
 $xyxy$
 $xyxyxyxyxy$

5

String $s = xyxxyxyxyxyxyxyxyxyxyxyxyxy$
 $xyxy$
 x
 $xyxyxyxyxy$

6

```
String s = xyxxxyxyxyyyxyxyxyxyxyxx
           xyxy
           x
           xy
           xxyyxyxyxx
```

```
String s = xyxxxyxyyyxyxyxyxyxyxx
           xyxy
           x
           xy
           xyxy
           xxyyxyxyxx
```

```
String s = xyxxxyxyxyxyxyxyxyxyxx
           xyxy
           x
           xy
           xyxy
           x
           xxyyxyxyxx
```

```
String s = xyxxxyxyxyyyxyxyxyxyxx
           xyxy
           x
           xy
           xyxy
           x
           xyxyxyxyxx
           xxyyxyxyxx
```

```
String s = xyxxxyxyxyxyxyxyxyxyxx
           xyxy
           x
           xy
           xyxy
           x
           xyxyxyxyxx
           x
           xyxyxyxyxx
```

```
String s = xyxxxyxyxyxyxyxyxyxyxx
           xyxy
           x
           xy
           xyxy
           x
           xyxyxyxyxx
           x
           xyx
           xxyyxyxyxx
```

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 x y x 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 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 x y x y x x

```

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 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 x y x y x x

```

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 y x y y
x
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 x y x y x x

```

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

```

Worst-case time $O(mn)$

16

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

```

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

```

Lots of wasted work

17

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 $O(m)$ states
- Run the DFA on the string $O(n)$
- Obvious construction method for DFA will be $O(m^2)$ but can be done in $O(m)$ time.
- Total $O(m+n)$ time

18

Building a DFA for the pattern

Pattern $p = xyxyxyxyxyxx$

19

Preprocessing the pattern

Pattern $p = xyxyxyxyxyxx$

20

Preprocessing the pattern

Pattern $p = xyxyxyxyxyxx$

21

Preprocessing the pattern

Pattern $p = xyxyxyxyxyxx$

22

Preprocessing the pattern

Pattern $p = xyxyxyxyxyxx$

23

Knuth-Morris-Pratt Algorithm

- Once the preprocessing is done there are only n steps on any string of size n
 - just follow your nose
- Obvious algorithm for doing preprocessing is $O(m^2)$ steps
 - still usually good since $m \ll n$
- Knuth-Morris-Pratt Algorithm can do the pre-processing in $O(m)$ steps
 - Total $O(m+n)$ time

24

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 with subset construction for the obvious NFA.)
- Typical text searches are based on finite automata designs
 - Perl builds this in as a first-class component of the programming language
 - `grep`

25