## String Matching

Find all occurrences of a pattern in a text

String Matching Problem:

Given text array T[1..n] and pattern array P[1..m] of characters from alphabet $\Sigma$, find all s such that T[s+1..s+m] = P[1..m], i.e., P occurs with shift s in T.

# Example

| r | o | w | | r | o | w | | r | o | w | | y | o | u | r | | b | o | a | t |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | P | = | y | o | | | | | | | |

$$s = 12$$
$$\Sigma = \{a, b, o, r, t, u, w, y\}$$

T: yoyoyoyo
P: yoyo

## String Matching

- Simple problem with many applications

  - text editing
  - pattern recognition

- Algorithms

  - Naive O((n-m+1)m) worst case
  - Rabin and Karp O((n-m+1)m) worst case, but better on average
  - Finite Automaton O(n+m| $\Sigma$ |)
  - Knuth-Morris-Pratt O(n+m)
  - Boyer and Moore O((n-m+1)m+| $\Sigma$ |) worst case, but better (best overall) in practice

---

## Naive String Matching

```
Naive(T, P)
  n = length(T)
  m = length(P)
  for s = 0 to n-m                          O(n-m+1)
     if P[1..m] = T[s+1..s+m]                  O(m)
     then print "Pattern occurs with shift" s
```

This algorithm takes O((n-m+1)m) time.

However, there is more information in a failed match:

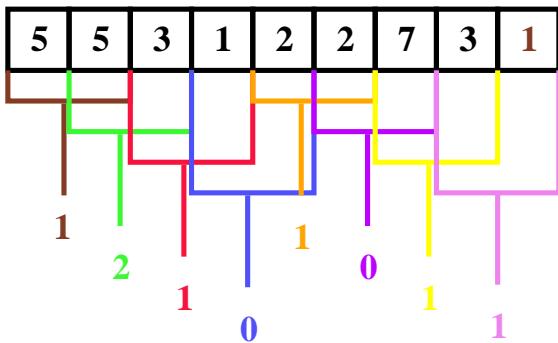| T: | a | a | a | a | b | a | a | ... |
|----|---|---|---|---|---|---|---|-----|
| P: | a | a | a | a | a |   |   |     |

$$s = s + m$$

No need to consider _____

---

## Rabin-Karp Algorithm

- Let characters be digits in radix-$|\Sigma|$ notation.

- Choose a prime number $q$ such that $|\Sigma| q$ fits within a computer word to speed computations.

- Algorithm:
  Compute (P mod q)
  Compute (T[s+1, .., s+m] mod q) for s = 0 .. n-m
  Test against P only those sequences in T having the same (mod q) value

- (T[s+1, .., s+m] mod q) can be incrementally computed by subtracting the high-order digit, shifting, adding the low-order bit, all in modulo q arithmetic.

## Example

$$\Sigma = \{0, 1, .., 9\}$$
$$P = 12, P \bmod 3 = 0$$
$$q = 3$$

| 5 | 5 | 3 | 1 | 2 | 2 | 7 | 3 | 1 |
|---|---|---|---|---|---|---|---|---|

1
  2
   1
    0
     1
      0
       1
        1

## Analysis

The Rabin-Karp algorithm takes $\Theta((n - m + 1)m)$ time in the worst case.
$O(n) + O(m(v + n/q))$ average case, v = #valid shifts
If q ≥ m and v = $O(1)$, then $O(n+m)$.
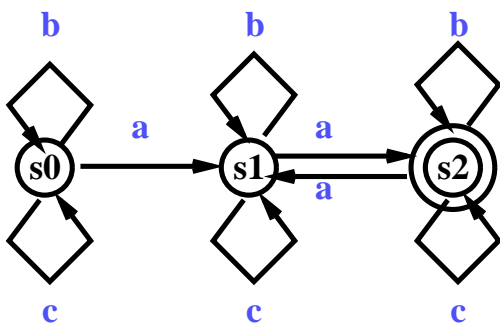
---

## Finite Automata

A finite automata M = $(Q, q_0, A, \Sigma, \delta)$, where

- Q = set of states $(s_i)$

- $q_0$ = start state $(s_0)$

- A = set of accepting states

- $\Sigma$ = input alphabet

- $\delta$ = transition function $Q\mathrm{x}\Sigma \to Q$

---

## Example

Here is a finite automaton accepting strings with an even number of "a"s.
$\Sigma = \{a, b, c\}$.

$$\delta(s_0, \text{a}) = s_1$$
$$\delta(s_0, \text{b}) = \delta(s_0, \text{c}) = s_0$$
$$\delta(s_1, \text{a}) = s_2$$
$$\delta(s_1, \text{b}) = \delta(s_1, \text{c}) = s_1$$
$$\delta(s_2, \text{a}) = s_1$$
$$\delta(s_2, \text{b}) = \delta(s_2, \text{c}) = s_2$$

$$A = \{s_2\}$$

Consider input string $w$. If $w$ ends at state $s \in A$, then the FA accepts $w$; otherwise, the FA rejects $w$.

Example: str = bccabaccaba

Accept

---

## String Matching FA

1. Compute FA accepting P (m+1 states)

2. Run FA with input string T, printing shift whenever accepting state is reached.

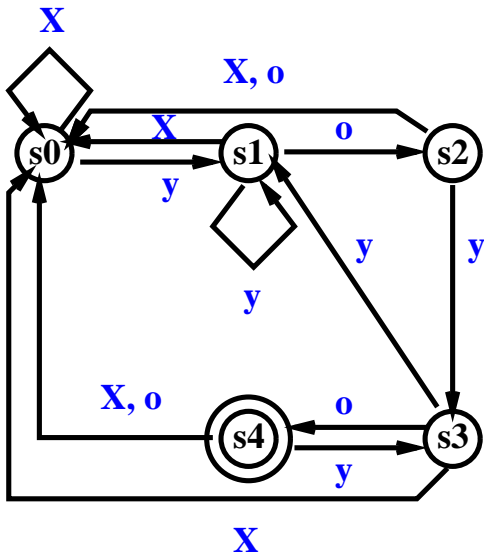## Example
$$P = \text{yoyo, m=4}$$
$$T = \text{spin your yoyo}$$
$$\Sigma = \{\text{i, n, o, p, r, s, u, y}\}$$
$$\text{Let X} = \Sigma - \{\text{y, o}\}$$

## Analysis

Computing $\delta$: $O(m|\Sigma|)$

FA-Matcher(T, $\delta$, m)          ; $O(n)$
    n = length(T)
    s = $s_0$
    for i = 1 to n
        s = $\delta$(s, T[i])
        if s = $s_m$
        then print "Pattern occurs with shift" (i-m)

        This algorithm takes $O(n + m|\Sigma|)$ time.

## Knuth-Morris-Pratt Algorithm

- Utilize a prefix array $\pi[1..m]$, where $\pi[q]$ contains information to compute $\delta(q, a)$ for $(a \in \Sigma)$, the pattern shift for a mismatch on P[q].

- $\pi$ requires only O(m) time (as opposed to $O(m|\Sigma|)$ for $\delta$).

---

## Prefix Array

# Example

| n | e | y | o | y | o | d | y | n | e | y | o |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | y | o | y | o | s |   |   |   |   |   |

$$s = 2$$

How far can we shift P over and be assured of catching all matches?

Since we have matched up to yoyo and yo is a suffix of yoyo, then we can shift over by 2 and start testing at P[3].

---

## Prefix Array

$\pi[q]$ answers the question:
  If we have matched P[1..q] in T, but P[q+1] does not match, then what is the longest prefix of P, P[1..k], that is a suffix of P[1..q]?
  We can then start matching again from P[k+1].

$$\pi[q] = \max\{k \mid k < q \text{ and } P[1..k] \text{ is a suffix of } P[1..q]\}$$

# Example

$$P = \begin{array}{ccccc} 1 & 2 & 3 & 4 & 5 \\ y & o & y & o & s \end{array}$$

$$\pi = \begin{array}{ccccc} 1 & 2 & 3 & 4 & 5 \\ 0 & 0 & 1 & 2 & 0 \end{array}$$

---

# Pseudocode

Compute-Prefix-Function(P)
    m = length(P)
    $\pi[1] = 0$                    ; k must be less than q
    k = 0
    for q = 2 to m                ; O(m) amortized
        while k > 0 and P[k+1] $\neq$ P[q]
            k = $\pi[k]$
        if P[k+1] = P[q]
        then k = k + 1            ; prefix increased by one
        $\pi[q]$ = k
    return $\pi$

---

# Pseudocode

KMP-Matcher(T, P)
    n = length(T)

```
m = length(P)
π = Compute-Prefix-Function(P)          ; O(m) amortized
q = 0
for i = 1 to n                          ; O(n) amortized
    while q > 0 and P[q+1] ≠ T[i] ; where do we move to in P?
        q = π[q]
    if P[q+1] = T[i]                     ; matches so far
    then q = q + 1
    if q = m
    then print "Pattern occurs with shift" (i-m)
        q = π[q]
```
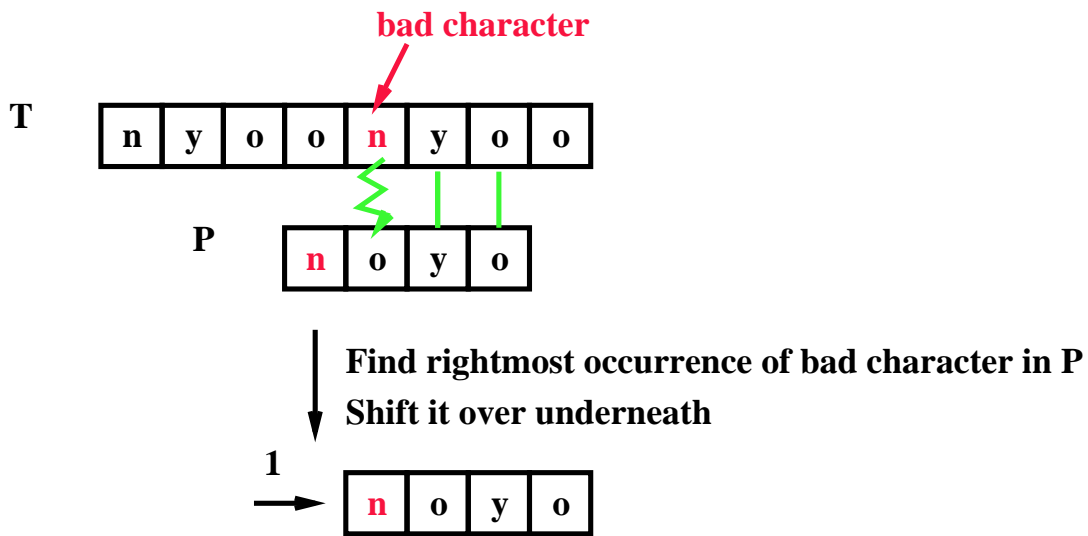
This algorithm takes _____ time

---

## Boyer-Moore Algorithm

- Most efficient (on average) when P is long and $\Sigma$ is large

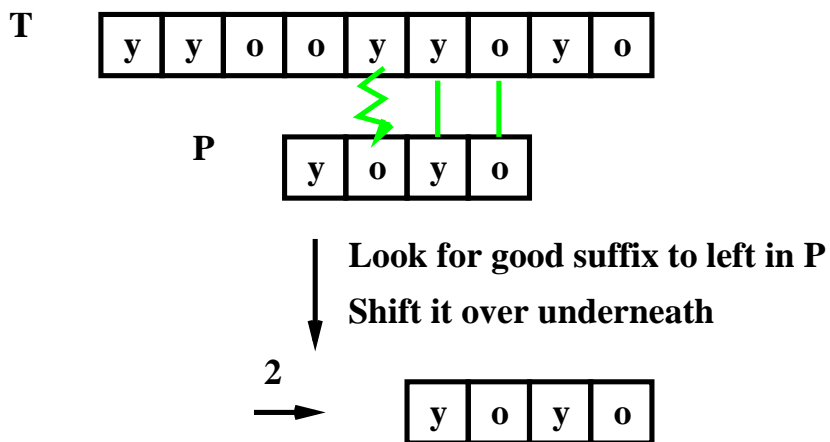- Matches pattern from right to left

- Utilizes two heuristics

# Bad Character Heuristic

Example

**T**

| n | y | o | o | n | y | o | o |
|---|---|---|---|---|---|---|---|

**P**

| n | o | y | o |
|---|---|---|---|

**Find rightmost occurrence of bad character in P**

**Shift it over underneath**

1 →

| n | o | y | o |
|---|---|---|---|

---

## Good Suffix Heuristic

Example

**T**

| y | y | o | o | y | y | o | y | o |
|---|---|---|---|---|---|---|---|---|

**P**

| y | o | y | o |
|---|---|---|---|

**Look for good suffix to left in P**

**Shift it over underneath**

2 →

| y | o | y | o |
|---|---|---|---|

---

# Information For Bad Character Heuristic

Compute-Last-Occurrence(P, m $\Sigma$)
    foreach a $\in \Sigma$
        $\lambda[a] = 0$
    for j = 1 to m
        $\lambda[P[j]] = j$
    return $\lambda$

$$\text{Running time: } O(|\Sigma| + m)$$

If mismatch at P[j] $\neq$ T[s+j], then shift (j - $\lambda$[T[s+j]]).

**Note:** Shift could be negative, in which case ignore the shift value and use Good Suffix shift which always has a positive value.

---

# Information for Good Suffix Heuristic

$$\gamma[j] = m - \max\{k \mid 0 \leq k < m \text{ and } P[j+1..m] \sqsupset P_k \text{ or } P_k \sqsupset P[j+1..m]\}$$
$$\sqsupset \text{ means } suffix \text{ (note: } x \sqsupset x)$$

If match j+1..m and P[j] $\neq$ T[s+j], shift right $\geq \gamma[j]$

# Examples

googoo

| 3 | 3 | 3 | 3 | 3 | 1 | 1 |
|---|---|---|---|---|---|---|

j = 0, $P_3 \sqsupset$ P[1..6]

googo

| 3 | 3 | 3 | 3 | 2 | 1 |
|---|---|---|---|---|---|

---

## Pseudocode

Compute-Good-Suffix(P, m)
    $\pi$ = Prefix(P)
    P' = reverse(P)
    $\pi$' = Prefix(P')
    for j = 0 to m                  ; O(m)
        $\gamma$[j] = m - $\pi$[m]
    for l = 1 to m
        j = m - $\pi$'[l]
        if $\gamma$[j] > l - $\pi$'[l]
        then $\gamma$[j] = l - $\pi$'[l]
    return $\gamma$

# Example

$$m = 4$$
$$P = yoyo, \pi = \underline{\qquad}$$
$$P' = oyoy, \pi' = \underline{\qquad}$$
$$\gamma = \underline{\qquad}$$
$$\gamma = \underline{\qquad}$$

---

## Boyer-Moore-Matcher

Boyer-Moore-Matcher(T, P, $\Sigma$)
    n = length(T)
    m = length(P)
    $\lambda$ = Compute-Last-Occurrence(P, m, $\Sigma$)     ; O($|\Sigma| + $ m)
    $\gamma$ = Compute-Good-Suffix(P, m)            ; O(m)
    s = 0

```
while s ≤ n-m                                    ; O(n-m+1)
    j = m
    while j > 0 and P[j] = T[s+j]                ; O(m)
        j = j - 1
    if j = 0
    then print "Pattern occurs with shift" s
        s = s + γ[0]
    else s = s + max(γ[j], j - λ[T[s+j]])
```

Close to naive
$O((n\text{-}m\text{+}1)m + \mid \Sigma \mid)$
Boyer-Moore-Matcher is actually best in practice

---

## Example

$$T = \text{soyoyo}$$
$$P = \text{yoyo}$$
$$\gamma = \underline{\hspace{2cm}}$$
$$\Sigma = \{o, s, y\}$$
$$\lambda = \underline{\hspace{1cm}}$$
$$\gamma = \underline{\hspace{3cm}}$$
$$\xrightarrow{2} \text{yoyo}$$
Match

---

## Applications