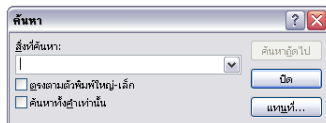# String Matching/ Pattern Matching algorithm

**Mr. Akarapon Watcharapalakorn**
**Doctor of Optometry (O.D.)**

---

# Syllabus

String Matching/Pattern Matching algorithm
- Brute Force algorithm (Naive)
- Rabin-Karp algorithm
- Finite Automaton algorithm
- Knutt-Morris-Pratt algorithm
- Boyer-Moore algorithm
- Horspool algorithm

---



---

# Exact String Matching

You are the fairest of your sex,
Let me be your hero;
I love you as one over $x$,
As $x$ approaches zero.
Positively.

text    pattern    you

Input

You are the fairest of your sex,
Let me be your hero;
I love you as one over $x$,
As $x$ approaches zero.
Positively.

Output

---

# Approximate String Match.

You are the fairest of your sex,
Let me be your hero;
I love you as one over $x$.
As $x$ approaches zero.
Positively.

text    pattern    heero

Input

You are the fairest of your sex,
Let me be your hero;
I love you as one over $x$,
As $x$ approaches zero.
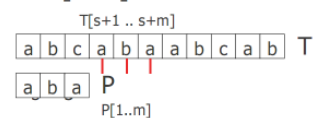Positively.

Output

---

# ปัญหา String Matching

- **Finite alphabet : $\Sigma$**
  - $\Sigma$ = {a, b, c, ..., z}
  - $\Sigma$ = {0, 1}
  - $\Sigma$ = {A,C,G,T}, ...
- **Text** : T[1..n]
- **Pattern** : P[1..m]

ต้องการหา
all valid shifts
first valid shift
last valid shift

$0 \le s \le n - m$

T[s+1 .. s+m]

| a | b | c | a | b | a | a | b | c | a | b | T
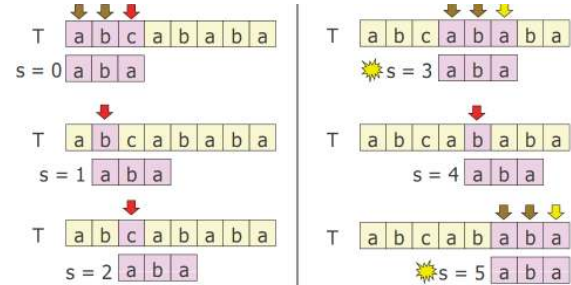
| a | b | a | P

P[1..m]

s เป็น <u>valid shift</u> เมื่อ P[1..m] = T[s+1 .. s+m]

1

## อัลกอริทึมการจับคู่สตริง

- Brute Force algorithm
- Deterministic Finite Automaton algorithm
- Rabin-Karp algorithm
- Shift Or algorithm
- Morris-Pratt algorithm
- Knuth-Morris-Pratt algorithm
- Simon algorithm
- Colussi algorithm
- Galil-Giancarlo algorithm
- Apostolico-Crochemore algorithm
- Boyer-Moore algorithm
- Turbo BM algorithm
- Apostolico-Giancarlo algorithm
- Reverse Colussi algorithm
- Horspool algorithm
- Quick Search algorithm
- Tuned Boyer-Moore algorithm
- Zhu-Takaoka algorithm
- Berry-Ravindran algorithm
- Smith algorithm
- Raita algorithm
- Reverse Factor algorithm
- Turbo Reverse Factor algorithm
- Forward Dawg Matching algorithm
- Backward Nondeterministic Dawg Matching algorithm
- Backward Oracle Matching algorithm
- Galil-Seiferas algorithm
- Two Way algorithm
- String Matching on Ordered Alphabets algorithm
- Optimal Mismatch algorithm
- Maximal Shift algorithm
- Skip Search algorithm
- KMP Skip Search algorithm

## Brute Force(Naive) Algo.



## Brute Force(Naive) Algo.

```
Naive-String-Matching( T[1..n], P[1..m] ) {
    for( s = 0 to n-m ) {          ← Θ(n−m)
        for( i = 1 to m ) {
            if (T[s+i] ≠ P[i]) break;   O(m)
        }
        if (i > m) print( s )
    }
}
```

$$O(nm)$$

worst case :    T = AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAZ
                P = AAAAZ

average # of comparisons : (random text and random pattern)
d = |Σ|

$$\frac{1-d^{-m}}{1-d^{-1}}(n-m+1) < 2(n-m+1)$$

## Random T&P : Avg #Cmps

| ตัวที่ 1 ไม่เหมือน | $\left(\dfrac{d-1}{d}\right)$ | $\times 1$ |
| ตัวแรกเหมือน ตัวที่ 2 ไม่เหมือน | $\left(\dfrac{1}{d}\right)^{1}\left(\dfrac{d-1}{d}\right)$ | $\times 2$ |
| สองตัวแรกเหมือน ตัวที่ 3 ไม่เหมือน | $\left(\dfrac{1}{d}\right)^{2}\left(\dfrac{d-1}{d}\right)$ | $\times 3$ |
| ... | | |
| m − 1 ตัวแรกเหมือน ตัวที่ m ไม่เหมือน | $\left(\dfrac{1}{d}\right)^{m-1}\left(\dfrac{d-1}{d}\right)$ | $\times m$ |
| เหมือนกันทั้ง m ตัว | $\left(\dfrac{1}{d}\right)^{m}$ | $\times m$ |

$$\sum_{k=1}^{m}\frac{k}{d^{k-1}}\left(\frac{d-1}{d}\right)$$
$$+$$
$$\frac{m}{d^{m}}$$

## Random T&P : Avg #Cmps

$$\frac{m}{d^{m}}+\sum_{k=1}^{m}\frac{k}{d^{k-1}}\left(\frac{d-1}{d}\right)=\frac{m}{d^{m}}+\left(1-\frac{1}{d}\right)\sum_{k=1}^{m}\frac{k}{d^{k-1}}=\,?$$

$$\sum_{k=0}^{m}x^{k}=\frac{1-x^{m+1}}{1-x}$$

$$\sum_{k=1}^{m}kx^{k-1}=\frac{-(1-x)(m+1)x^{m}+(1-x^{m+1})}{(1-x)^{2}}$$

$$=\frac{1-(m+1)x^{m}+mx^{m+1}}{(1-x)^{2}}$$

$$=\frac{1-x^{m}+mx^{m}(x-1)}{(1-x)^{2}}$$

## Random T&P : Avg #Cmps

$$\frac{m}{d^{m}}+\left(1-\frac{1}{d}\right)\sum_{k=1}^{m}\frac{k}{d^{k-1}}=\frac{1-d^{-m}}{1-d^{-1}}$$

$$\sum_{k=1}^{m}kx^{k-1}=\frac{1-x^{m}+mx^{m}(x-1)}{(1-x)^{2}}$$

$$mx^{m}+(1-x)\sum_{k=1}^{m}kx^{k-1}=mx^{m}+(1-x)\left(\frac{1-x^{m}+mx^{m}(x-1)}{(1-x)^{2}}\right)$$

$$=mx^{m}+\left(\frac{1-x^{m}+mx^{m}(x-1)}{(1-x)}\right)$$

$$=\frac{1-x^{m}}{1-x}\quad\longrightarrow\quad \text{ให้ } x=d^{-1}$$

## Random T&P : Avg #Cmps

```
Naive-String-Matching( T[1..n], P[1..m] ) {
  for( s = 0 to n-m ) {
    for( i = 1 to m ) {
      if (T[s+i] ≠ P[i]) break;
    }
    if (i > m) print( s )
  }
}
```
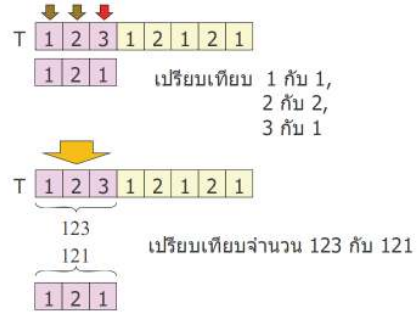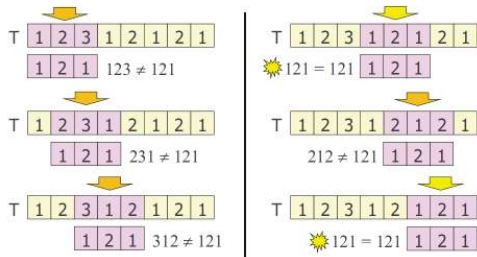
$$\frac{1-d^{-m}}{1-d^{-1}}$$

$$\frac{1-d^{-m}}{1-d^{-1}}(n-m+1) < 2(n-m+1)$$

$$\frac{1-d^{-m}}{1-d^{-1}} < \frac{1}{1-d^{-1}}, \ d \geq 2$$

$$\leq \frac{1}{1-2^{-1}} = 2$$

## Rabin-Karp Matcher



เปรียบเทียบ  1 กับ 1,
2 กับ 2,
3 กับ 1

เปรียบเทียบจำนวน 123 กับ 121

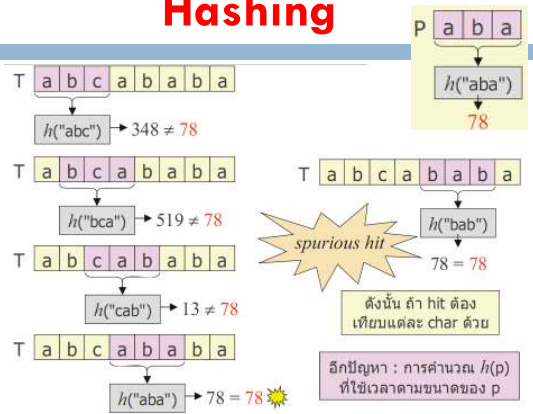## Rabin-Karp Algorithm



ถ้า pattern เป็นตัวอักษร  หรือ
ถ้าเป็นตัวเลข แต่มีขนาดใหญ่เกิน int เกิน long จะทำอย่างไร ?

## Hashing



spurious hit

ดังนั้น ถ้า hit ต้อง
เทียบแต่ละ char ด้วย

อีกปัญหา : การคำนวณ $h(p)$
ที่ใช้เวลาตามขนาดของ p

## การคำนวณค่า h: incremental

❖ $h(s)$ = (มอง s เป็นจำนวน ในระบบเลขฐาน)
❖ $h(\text{"abc"}) = (1 \times 37^2 + 2 \times 37^1 + 3 \times 37^0) = 1446$
❖ $h(\text{"bce"}) = (2 \times 37^2 + 3 \times 37^1 + 5 \times 37^0) = 2853$
$$= 37 \times (1 \times 37^2 + 2 \times 37^1 + 3 \times 37^0) + 5 \times 37^0$$
$$= 37 \times (h(\text{"abc"}) - 1 \times 37^2) + 5 \times 37^0$$

เลื่อนตำแหน่ง "bc"
ไปทางซ้าย    ลบ "a"    เพิ่ม "e"

$$h(\text{"...e"}) = 37 \times (h(\text{"a..."}) - 1 \times 37^{m-1}) + 5 \times 37^0$$

## Rabin-Karp Algorithm

```
Rabin-Karp-Matching( T[1..n], P[1..m], d ) {
  ht = 0; hp = 0                              d = |Σ|
  for ( i = 1 to m ) {
    ht = d*ht + T[i]
    hp = d*hp + P[i]
  }
  for ( s = 0 to n - m ) {
    if (ht == hp) {
      for ( i = 1 to m )
        if (T[s+i] ≠ P[i]) break
      if (i > m) print( s )
    }
    if (s < n-m) {
      ht = d*(ht - T[s+1]*d^(m-1)) + T[s+m+1]
    }
  }
}       h("...e") = 37×( h("a...") − 1×37^(m−1) ) + 5×37^0
```

ถ้า m มีค่ามาก ขนาด
ของ ht และ hp ก็ใหญ่

# Rabin-Karp Algorithm

```
Rabin-Karp-Matching( T[1..n], P[1..m], d, q ) {
  ht = 0; hp = 0; dm1 = d^(m-1) % q
  for ( i = 1 to m ) {
    ht = (d*ht + T[i]) % q
    hp = (d*hp + P[i]) % q
  }
  for ( s = 0 to n - m ) {
    if (ht == hp) {
      for ( i = 1 to m )
        if (T[s+i] ≠ P[i]) break
      if (i > m) print( s )
    }
    if (s < n-m) {
      ht = (d*(ht - T[s+1]*dm1) + T[s+m+1]) % q
    }
  }
}
```
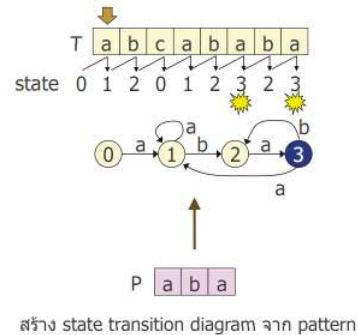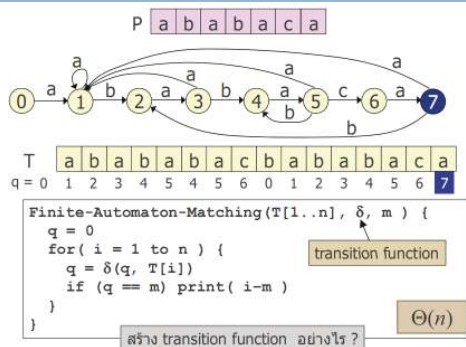
worst case $\Theta(nm)$

avg case $O(n + m)$

$h("...e") = 37 \times ( h("a...") - 1 \times 37^{m-1} ) + 5 \times 37^0$

---

# Finite Automaton Matcher



สร้าง state transition diagram จาก pattern

---

# Finite Automaton Matcher



```
Finite-Automaton-Matching(T[1..n], δ, m ) {
  q = 0
  for( i = 1 to n ) {
    q = δ(q, T[i])
    if (q == m) print( i-m )
  }
}
```

transition function

$\Theta(n)$

สร้าง transition function อย่างไร ?

---

# Transition Function



$\Sigma = \{a, b, c\}$

ตารางมีขนาด $(m+1)|\Sigma|$

$\delta$(state, input)

---

# Transition Function



$\delta(q, x) = \max\{ k : (P[1..q]+x).endsWith(P[1..k]) \}$

---

# ตัวอย่าง Transition Function



$\delta(q, x) = \max\{ k : (P[1..q]+x).endsWith(P[1..k]) \}$

$\delta(5,"b") = 4$

$P[1..5]+"b" = ababab$

สร้างแบบตรงไปตรงมาใช้เวลา $O(m^3|\Sigma|)$

มีวิธีสร้างที่ใช้เวลา $O(m|\Sigma|)$

รวมการค้นด้วย ใช้เวลา $O(n + m|\Sigma|)$

# Knutt-Morris-Pratt Matcher

- ❖ ของเดิม (transition function)
  - ❖ $\delta(q, x)$ = max{ k : (P[1..q]+x).endsWith(P[1..k]) }
  - ❖ ตารางมีขนาด $(m+1)|\Sigma|$
  - ❖ ใช้เวลาสร้าง $O(m|\Sigma|)$   | 1,2,3, | ..., q | x |
  - | 1,2,..., k |
- ❖ ของใหม่ (prefix function)
  - ❖ $\pi(q)$ = max{ k : k<q and (P[1..q]).endsWith(P[1..k]) }
  - ❖ match มาได้ q ตัว แต่ mismatch ตัวที่ q+1 ให้ทำต่อที่ $P[\pi(q)+1]$
  - ❖ ตารางมีขนาด m   | 1,2,3, | ... , q |
  - ❖ ใช้เวลาสร้าง $\Theta(m)$   | 1,2,...,k |
- ❖ สร้าง $\pi$ และค้น ใช้เวลารวม $\Theta(n+m)$

# การใช้ Prefix Function



# KMP Matcher

```
KMP-Matcher( T[1..n], P[1..m] ) {
    π = Compute-Prefix-Function(P)  // Θ(m)
    q = 0
    for ( i = 1 to n ) {
        while (q > 0 AND P[q+1] ≠ T[i]) {
            q = π[q]
        }
        if (P[q+1] == T[i]) q++
        if (q == m) {
            print( i - m )      Θ(n + m)
            q = π[q]
        }
    }
}
```



$\pi[4] = 2$

## $\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).endsWith(P[1..k]) \}$



## $\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).endsWith(P[1..k]) \}$



```
Compute-Prefix-Function( P[1..m] ) {
    π[1] = 0;  k = 0
    for ( q = 2 to m ) {

        ...

        if (P[k+1] == P[q]) k++
        π[q] = k
    }
    return π
```

## $\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).endsWith(P[1..k]) \}$



```
Compute-Prefix-Function( P[1..m] ) {
    π[1] = 0;  k = 0
    for ( q = 2 to m ) {
        while (k > 0 and P[k+1] ≠ P[q]) {
            k = π[k]
        }
        if (P[k+1] == P[q]) k++
        π[q] = k
    }
    return π
```

## Slide 1

$\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).\text{endsWith}(P[1..k]) \}$

| a | t | c | a | c | a | t | c | a | t | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|

$\pi[\ 7\ ] = 2$

$\pi[\ 8\ ] = 3$

$\pi[\ 9\ ] = 4$

```
Compute-Prefix-Function( P[1..m] ) {
    π[1] = 0;  k = 0
    for ( q = 2 to m ) {
        while (k > 0 and P[k+1] ≠ P[q]) {
            k = π[k]
        }
        if (P[k+1] == P[q]) k++
        π[q] = k
    }
    return π
```

## Slide 2

$\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).\text{endsWith}(P[1..k]) \}$

| a | t | c | a | c | a | t | c | a | t | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|

$\pi[\ 9\ ] = 4 \ \rightarrow\ P[1..4] = P[6..9] \rightarrow$ ทดสอบ $P[4+1] = P[10]$ ?

$\pi[\ 4\ ] = 1 \ \rightarrow\ P[1..1] = P[4..4] \rightarrow$ ทดสอบ $P[1+1] = P[10]$ ?

$\pi[\ 10\ ] = 2$

```
Compute-Prefix-Function( P[1..m] ) {
    π[1] = 0;  k = 0
    for ( q = 2 to m ) {
        while (k > 0 and P[k+1] ≠ P[q]) {
            k = π[k]
        }
        if (P[k+1] == P[q]) k++
        π[q] = k
    }
    return π
```

## Slide 3

$\pi(q) = \max\{ k : k<q \text{ and } (P[1..q]).\text{endsWith}(P[1..k]) \}$

| a | t | c | a | c | a | t | c | a | t | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|

$\pi[\ 11\ ] = 3$

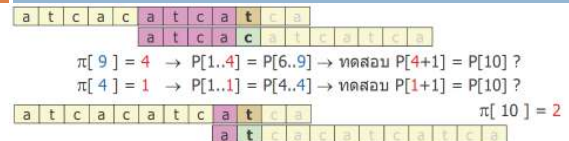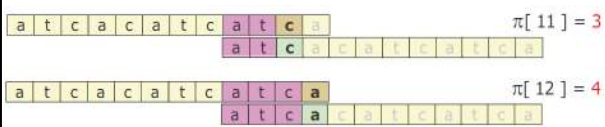$\pi[\ 12\ ] = 4$

```
Compute-Prefix-Function( P[1..m] ) {
    π[1] = 0;  k = 0
    for ( q = 2 to m ) {
        while (k > 0 and P[k+1] ≠ P[q]) {
            k = π[k]
        }
        if (P[k+1] == P[q]) k++
        π[q] = k
    }
    return π
```

$\Theta(m)$

## Slide 4
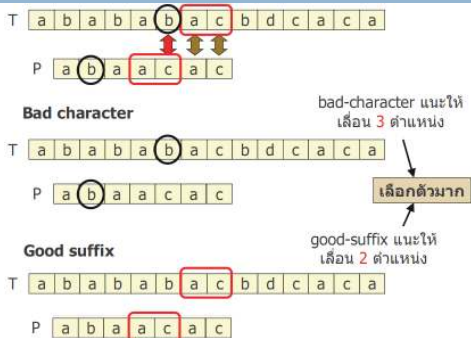
# Boyer-Moore Matcher

**Knuth-Morris-Pratt (left-to-right)**
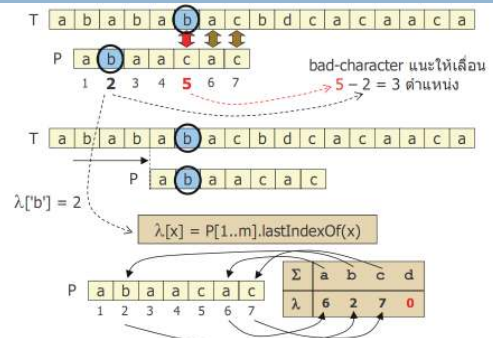
A STRING SEARCHING EXAMPLE CONSISTING OF ...

... STING

**Boyer Moore (right-to-left)**
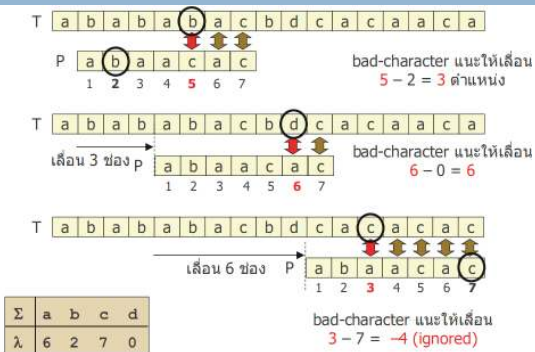
A STRING SEARCHING EXAMPLE CONSISTING OF ...

STING

## Slide 5

# Boyer-Moore: two heuristics

| T | a | b | a | b | a | b | a | c | b | d | c | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|

**Bad character**

| T | a | b | a | b | a | b | a | c | b | d | c | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|

bad-character แนะให้
เลื่อน 3 ตำแหน่ง

เลือกตัวมาก

**Good suffix**

| T | a | b | a | b | a | b | a | c | b | d | c | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|

good-suffix แนะให้
เลื่อน 2 ตำแหน่ง

## Slide 6

# Bad-Character Heuristics

| T | a | b | a | b | a | b | a | c | b | d | c | a | c | a | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

bad-character แนะให้เลื่อน

$5 - 2 = 3$ ตำแหน่ง

| T | a | b | a | b | a | b | a | c | b | d | c | a | c | a | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|

$\lambda['b'] = 2$

$\lambda[x] = P[1..m].\text{lastIndexOf}(x)$

| P | a | b | a | a | c | a | c |
|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| Σ | a | b | c | d |
|---|---|---|---|---|
| λ | 6 | 2 | 7 | 0 |

# Bad-Character Heuristics

T: a b a b a b a c b d c a c a a c a
P: a b a a c a c
1 2 3 4 5 6 7

bad-character แนะให้เลื่อน
5 − 2 = 3 ตำแหน่ง

T: a b a b a b a c b d c a c a a c a
เลื่อน 3 ช่อง   P: a b a a c a c
1 2 3 4 5 6 7

bad-character แนะให้เลื่อน
6 − 0 = 6

T: a b a b a b a c b d c a c a c a c
เลื่อน 6 ช่อง   P: a b a a c a c
1 2 3 4 5 6 7

bad-character แนะให้เลื่อน
3 − 7 = −4 (ignored)

| Σ | a | b | c | d |
|---|---|---|---|---|
| λ | 6 | 2 | 7 | 0 |

---

# Good-Suffix Heuristics

T: .. .. .. .. .. b a c .. .. .. .. .. .. .. ..
P: a a c a d a c
1 2 3 4 5 6 7

$\gamma[\ 5\ ] = 4$

mismatch ตำแหน่งที่ 5
good-suffix แนะให้เลื่อน
4 ตำแหน่ง

T: .. .. .. .. .. b a c .. .. .. .. .. ..
P: a a c a d a c
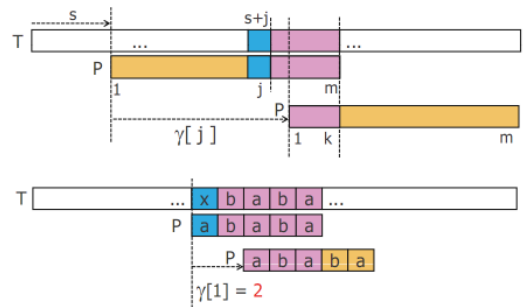
---

# Good-Suffix Heuristics

$\gamma[\,j\,] = m - \max \{\ k : 0 \leq k < m \text{ and } P[1..k].\text{endsWith}(P[j+1...m])\ \}$

P: a b a b a
P: a b a b a
$\gamma[4] = 2$

---

# Good-Suffix Heuristics

$\gamma[\,j\,] = m - \max \{\ k : 0 \leq k < m \text{ and } P[j+1...m].\text{endsWith}(P[1..k])\ \}$

P: a b a b a
P: a b a b a
$\gamma[1] = 2$

---

$\gamma[\,j\,] = m - \max \{\ k : 0 \leq k < m \text{ and } P[j+1...m] \sim P[1..k]\ \}$

P: A G A G G A G

$\gamma[7] = 1$

j=6   A G A G G A G   $\gamma[6] = 7 - 5 = 2$

j=5   A G A G G A G   $\gamma[5] = 7 - 4 = 3$

j=4   A G A G G A G   $\gamma[4] = 7 - 4 = 3$

j=3   A G A G G A G   $\gamma[3] = 7 - 2 = 5$

j=2   A G A G G A G   $\gamma[2] = 7 - 2 = 5$

j=1   A G A G G A G   $\gamma[1] = 7 - 2 = 5$

j=0   A G A G G A G   $\gamma[0] = 7 - 2 = 5$

---

# Boyer-Moore Matcher

```
Boyer-Moore-Matcher(T[1..n],P[1..m], Σ)
   λ = LAST-OCCURRENCE(P, m, Σ)   // O(|Σ|+m)
   γ = GOOD-SUFFIX(P, m)          // O(m)
   s = 0
   while( s ≤ n - m ) {
      j = m
      while( j > 0 and P[j] = T[s+j]) j--
      if (j==0) {
         print( s )
         s = s + γ[0]
      } else {
         s = s + max(γ[j], j - λ[T[s+j]])
      }
   }
}
```

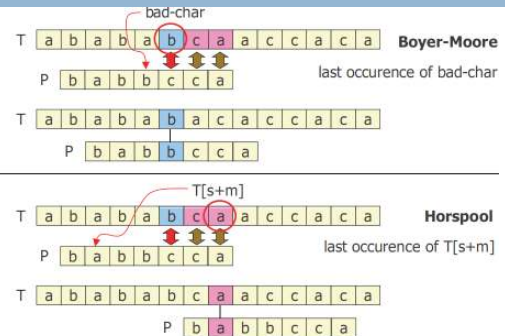worst : $O((n-m+1)m + |\Sigma|)$     best : $O(n/m + m + |\Sigma|)$
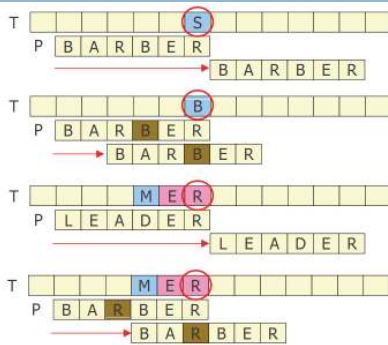
## Boyer-Moore Matcher

```
Boyer-Moore-Matcher(T[1..n],P[1..m], Σ)
    λ = LAST-OCCURRENCE(P, m, Σ)
    γ = GOOD-SUFFIX(P, m)
    s = 0
    while( s ≤ n - m ) {
        j = m
        while( j > 0 and P[j] = T[s+j]) j--
        if (j==0) {
            print( s )
            s = s + γ[0] 1
        } else {
            s = s + max(γ[j], j - λ[T[s+j]]) 1
        }
    }
}
```
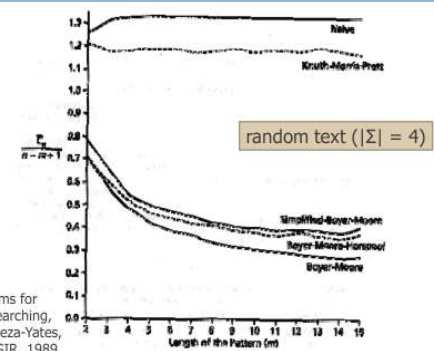
## Horspool Matcher



Boyer-Moore
last occurence of bad-char

Horspool
last occurence of T[s+m]

## Horspool Matcher



## ประสิทธิภาพการทำงาน



random text (|Σ| = 4)

Algorithms for string searching, R. A. Baeza-Yates, ACM SIGIR. 1989

## ประสิทธิภาพการทำงาน



English text

Algorithms for string searching, R. A. Baeza-Yates, ACM SIGIR, 1989

## THE END