# String Matching Algorithms

Dae-Ki Kang

# String Matching
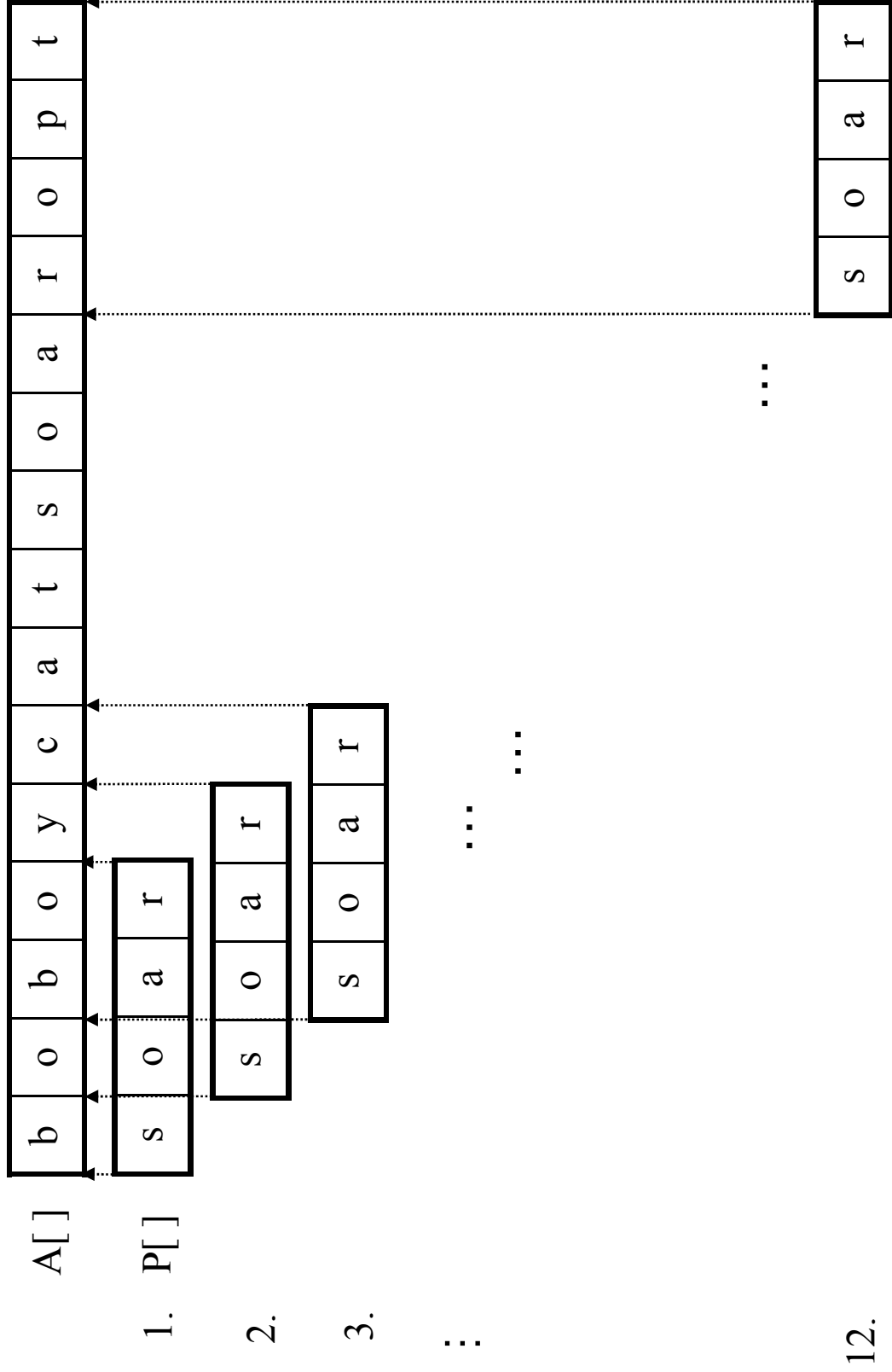
- Input
  - A[1...*n*]: Text string
  - P[1...*m*]: Pattern string
  - *m* << *n*
- Goal
  - Determine if A[1...*n*] includes P[1...m]

# Naïve Matching or Stupid Matching

naïveMatching(A[ ], P[ ])

{

△ *n*: length of A[ ] , *m*: length of P[ ]

**for** *i* ← 1 **to** *n-m*+1{

if (P[1...*m*] = A[*i*...*i+m*-1]) **then**

**There is a match at** A[*i*];

}

}

✔ Running time:  $O(mn)$

# How Naïve Matching Works?

A[] | b | o | b | o | y | c | a | t | s | o | a | r | o | p | t

1. P[] | s | o | a | r
2. | s | o | a | r
3. | s | o | a | r

…

12. | s | o | a | r

# When Naïve Matching is inefficient?

# Rabin-Karp Algorithm

- Change a pattern into a number and change a compared portion of a text string into a number
- Now, compare two numbers (instead of comparing two strings)
- String → Number
  - Depends on the cardinality of a character set $\sum$
  - Example: $\sum$ = {a, b, c, d, e}
    - $|\sum|$ = 5
    - Convert a, b, c, d, e into 0, 1, 2, 3, 4
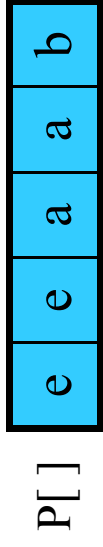    - String "cad" will be $2*5^2+0*5^1+3*5^0 = 28$

# Overhead of the Conversion

- Conversion time of A[i···i+m−1]

  – $a_i = A[i+m-1] + d(A[i+m-2] + d(A[i+m-3] + d(\cdots + d(A[i]))\cdots)$

  – $\Theta(m)$ time complexity

  – Total matching of A[1···n] takes $\Theta(mn)$

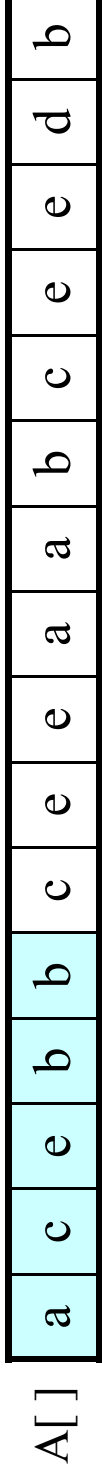  – No better than naïve matching

- Horner's rule

- Regardless of $m$, calculate the number as follows:

  – $a_i = d(a_{i-1} - d^{m-1}A[i-1]) + A[i+m-1]$

  – We calculate $d^{m-1}$ once

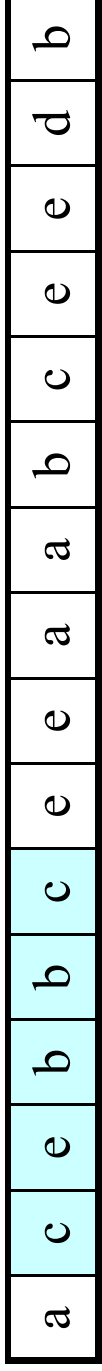  – Two multiplications and two additions

# One matching example

P[]  | e | e | a | a | b |

$p = 4*5^4 + 4*5^3 + 0*5^2 + 0*5^1 + 1 = 3001$

A[]  | a | c | e | b | c | e | e | a | a | b | c | e | e | d | b |

$a_1=0*5^4+2*5^3+4*5^2+1*5^1+1 = 356$

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_2= 5(a_1-0*5^4)+2 = 1782$

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_3=5(a_2-2*5^4)+4 = 2664$

...

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_7=5(a_6-2*5^4)+1 = 3001$

...

# Basic Rabin Karp Algorithm

basicRabinKarp(A, P, $d$, $q$)
{
  ▷ $n$ : length of A[ ], $m$ : length of P[ ]
  $p \leftarrow 0$; $a_1 \leftarrow 0$;
  for $i \leftarrow 1$ to $m$ {    ▷ calculate $a_1$
    $p \leftarrow dp + \text{P}[i]$;
    $a_1 \leftarrow da_1 + \text{A}[i]$;
  }
  for $i \leftarrow 1$ to $n-m+1$ {
    if ($i \neq 1$) then $a_i \leftarrow d(a_{i-1} - d^{m-1}\text{A}[i{-}1]) + \text{A}[i{+}\,m{-}1]$;
    if ($p = a_i$) then matched at A[$i$];
  }
}

✓ Total running time: $\Theta(n)$

# Problem of basic Rabin Karp Algorithm

- Depending on $|\Sigma|$ and $m$, $a_i$ can be big
  - bigger than CPU register size
  - and causes overflow

- Solution
  - Limit the size of $a_i$ using modulo operation
  - $a_i = d(a_{i-1} - d^{m-1}A[i-1]) + A[i+m-1]$ ➔
  - $b_i = (d(b_{i-1} - (d^{m-1} \bmod q) A[i-1]) + A[i+m-1]) \bmod q$
  - q is a big prime number, but $dq$ should be fit in a register

# Matching using modulo

P[ ]  | e | e | a | a | b |

$p = (4*5^4 + 4*5^3 + 0*5^2 + 0*5^1 + 1) \bmod 113 = 63$

A[ ]  | a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_1 = (0*5^4 + 2*5^3 + 4*5^2 + 1*5^1 + 1) \bmod 113 = 17$

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_2 = (5(a_1 - 0*(5^4 \bmod 113)) + 2) \bmod 113 = 87$

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_3 = (5(a_2 - 2*(5^4 \bmod 113)) + 4) \bmod 113 = 65$

...

| a | c | e | b | b | c | e | e | a | a | b | c | e | e | d | b |

$a_7 = (5(a_6 - 2*(5^4 \bmod 113)) + 1) \bmod 113 = 63$
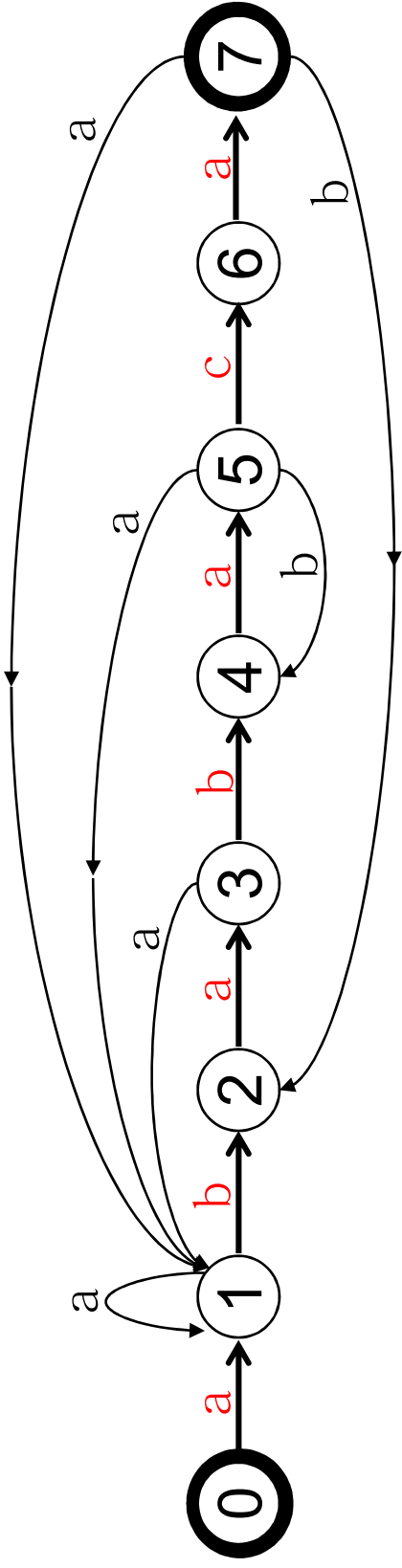
...

# Rabin–Karp Algorithm

RabinKarp(A, P, $d$, $q$)
{
$\triangleright$ $n$ : size of A[ ], $m$ : size of P[ ]
$p \leftarrow 0$; $b_1 \leftarrow 0$;
for $i \leftarrow 1$ to $m$ {                                      $\triangleright$ $b_1$
    $p \leftarrow (dp + P[i]) \bmod q$;
    $b_1 \leftarrow (db_1 + A[i]) \bmod q$;
}
$h \leftarrow d^{m-1} \bmod q$;
for $i \leftarrow 1$ to $n-m+1$ {
    if $(i \neq 1)$ then $b_i \leftarrow (d(b_{i-1} - hA[i-1]) + A[i+m-1]) \bmod q$;
    if $(p = b_i)$ then
        if $(P[1 \cdots m] = A[i \cdots i+m-1])$ then
            there is a match at A[$i$];
}
}

✓ **Average running time:** $\Theta(n)$

# Automata Based Matching

- Finite Automata
  - Finite symbols, states and transition
  - Five tuple: $(Q, q0, A, \Sigma, \delta)$
    - Q : states
    - Q0 : start state
    - A : accepted states
    - $\Sigma$ : input alphabet
    - $\delta$ : state transition diagram
- States represent a snapshot of matching process

# Automata that checks ababaca



S: dvganbbactababaa**ababaca**bababacaagbk...

# Implementation of Automata

| State / Input symbol | a | b | c | d | e | ... | z |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 0 | ... | 0 |
| 1 | 1 | 2 | 0 | 0 | 0 | ... | 0 |
| 2 | 3 | 0 | 0 | 0 | 0 | ... | 0 |
| 3 | 1 | 4 | 0 | 0 | 0 | ... | 0 |
| 4 | 5 | 0 | 0 | 0 | 0 | ... | 0 |
| 5 | 1 | 4 | 6 | 0 | 0 | ... | 0 |
| 6 | 7 | 0 | 0 | 0 | 0 | ... | 0 |
| 7 | 1 | 2 | 0 | 0 | 0 | ... | 0 |

| State / Input symbol | a | b | c | Else |
|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 |
| 1 | 1 | 2 | 0 | 0 |
| 2 | 3 | 0 | 0 | 0 |
| 3 | 1 | 4 | 0 | 0 |
| 4 | 5 | 0 | 0 | 0 |
| 5 | 1 | 4 | 6 | 0 |
| 6 | 7 | 0 | 0 | 0 |
| 7 | 1 | 2 | 0 | 0 |

# Algorithm that checks matching using automata

FA−Matcher (A, $\delta$ , $f$ )
▷ $f$ : accepted state
{
  ▷ $n$: length of A[ ]
  $q \leftarrow 0$;
  for $i \leftarrow 1$ to $n$ {
    $q \leftarrow \delta(q, A[i])$;
    if ($q = f$) then we found a match at A[$i-m+1$];
  }
}

✓ Total running time: $\Theta(n + |\Sigma|m)$
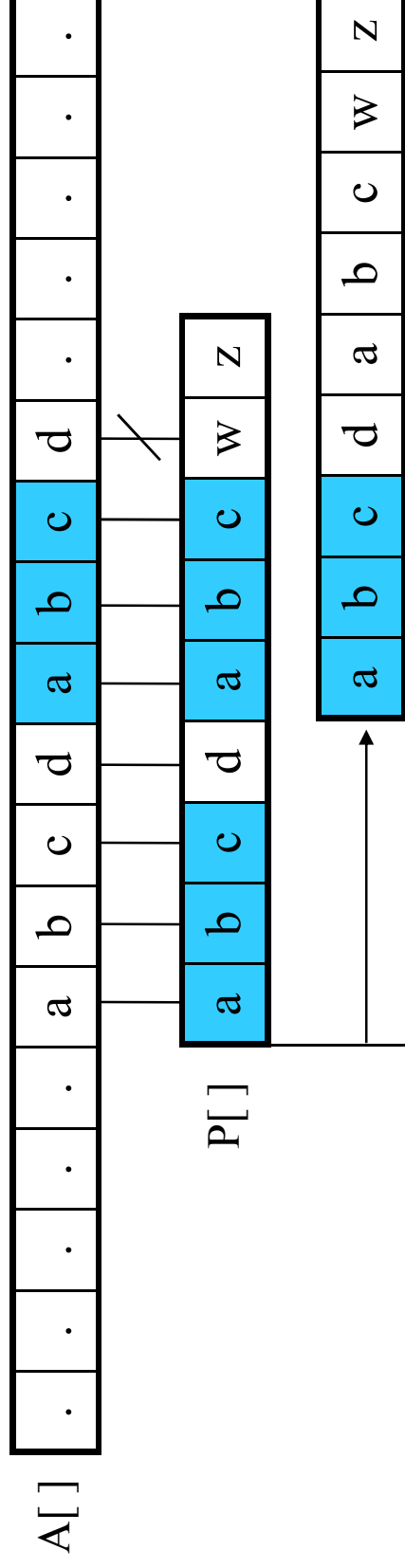
# Compute Transition Function

Compute_Transition_Function(P, $\Sigma$)
{
m ← |P|; $\triangle$ length of pattern
for q ← 0 to m do
    for each character a ∈ $\Sigma$ do
        k ← min(m+ 1,q+ 2);
        repeat k ← k − 1
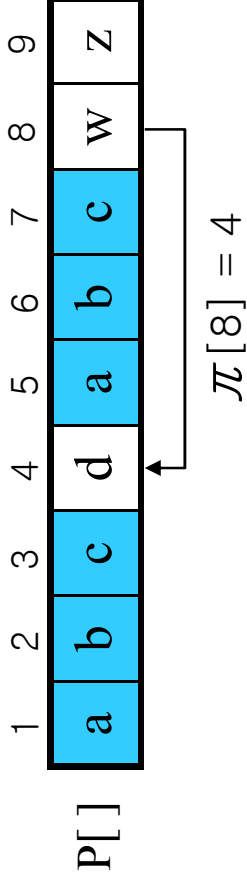        until P[1..k] is a suffix of P[1..q]a;
        $\delta(q, a)$ ← k;
}

✓Running time: $\Theta(|\Sigma|m)$

# KMP(Knuth-Morris-Pratt) Algorithm

- Similar to automata based matching

- Commonalities
  - Prepare states for mismatch
  - Simpler than automata based match

A[] 

| . | . | . | a | b | c | d | a | b | c | d | . | . | . | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

P[]

| a | b | c | d | a | b | c | w | z |
|---|---|---|---|---|---|---|---|---|

| a | b | c | d | a | b | c | w | z |
|---|---|---|---|---|---|---|---|---|

# Prepare return points for each mismatch

P[ ]

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|
| a | b | c | d | a | b | c | w | z |

$\pi[8] = 4$

Matched upto "abcdabc", failed at "w"
"abc" at 1,2,3 and "abc" at 5,6,7 are the same.
So, we compare mismatched text char with P[4]

For each index of pattern,
we prepare the return point.

$\pi[\ ]$

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|----|
| 0 | 1 | 1 | 1 | 1 | 2 | 3 | 4 | 1 | 1  |

| a | b | c | d | a | b | c | w | z |
|---|---|---|---|---|---|---|---|---|

# KMP Algorithm

KMP(A[ ], P[ ])
{
  preprocessing(P);
  $i \leftarrow 1$;  △ index pointer of text
  $j \leftarrow 1$;  △ index pointer of pattern
  △ $n$: size of A[ ], $m$: size of P[ ]
  while ($i \leq n$) {
      if ($j = 0$ or A[$i$] = P[$j$])
          then { $i++$;  $j++$; }
          else  $j \leftarrow \pi$[$j$];
      if ($j = m+1$) then {
          There is a match at A[$i-m$];
          $j \leftarrow \pi$[$j$];
      }
  }
}

✓Running time: $\Theta(n)$

# Preprocessing

```
preprocessing(P)
{
  m ← |P|;  ▷ length of pattern
  π[1] ← 0;
  k ← 0;
  for q ← 2 to m do
      while (k > 0) and (P[k+ 1] ≠ P[q]) do
          k ← π[k];
      if (P[k+ 1] = P[q]) then k ← k+ 1;
      π[q] ← k;
  return π;
}
```

✓Running time: $\Theta(m)$