# Exact String Matching: KMP & Boyer-Moore Algorithms
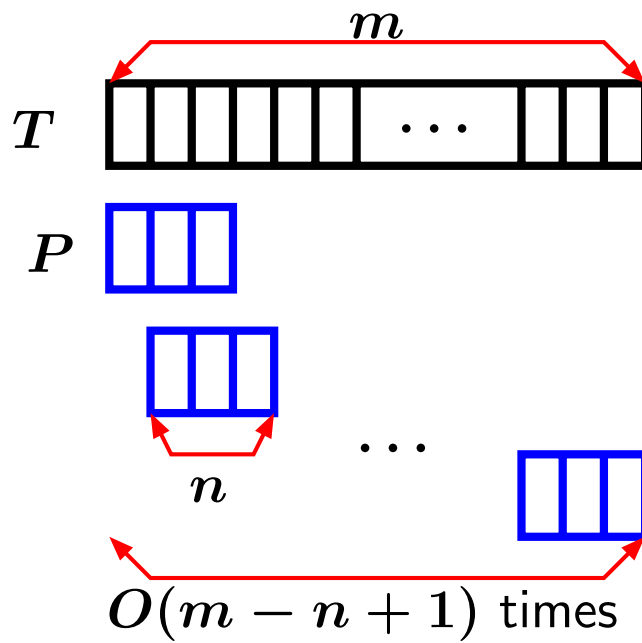
## R.C.T. Lee and Chin Lung Lu

CS 5313 Algorithms for Molecular Biology

# Exact string matching
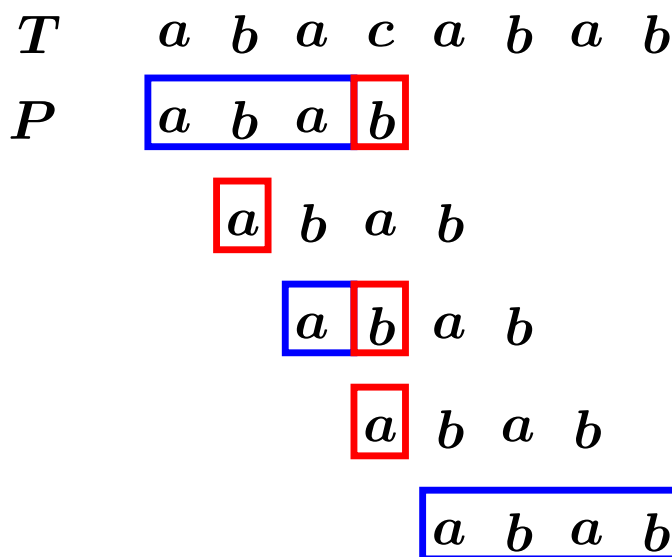
- Given two strings $T$ and $P$, where $|T| = m$ and $|P| = n$, find all the occurrences of $P$ in $T$?

- (**Simplified version**) Given two strings $T$ and $P$, if $P$ occurs in $T$?

- Brute-force method: $O(mn)$

- KMP algorithm (Knuth, Morris and Pratt, 1977): $\mathcal{O}(m+n)$

- Boyer-Moore algorithm (Boyer and Moore, 1977): $\mathcal{O}(m+n)$

# Brute-force method



Time complexity of brute-force method: $\mathcal{O}(mn)$

---

# Brute-force method

$$T \quad a \quad b \quad a \quad c \quad a \quad b \quad a \quad b$$
$$P \quad \boxed{a \quad b \quad a \quad \boxed{b}}$$
$$\boxed{a} \quad b \quad a \quad b$$
$$\boxed{a \quad \boxed{b}} \quad a \quad b$$
$$\boxed{a} \quad b \quad a \quad b$$
$$\boxed{a \quad b \quad a \quad b}$$

Simple, but not efficient because it cannot avoid rescanning $T$

# How to speedup the brute-force method?

- Shift $P$ by more than one place when a mismatch occurs without missing any occurrence of $P$ in $T$

1. KMP algorithm
   - Proposed by Knuth, Morris and Pratt at 1977
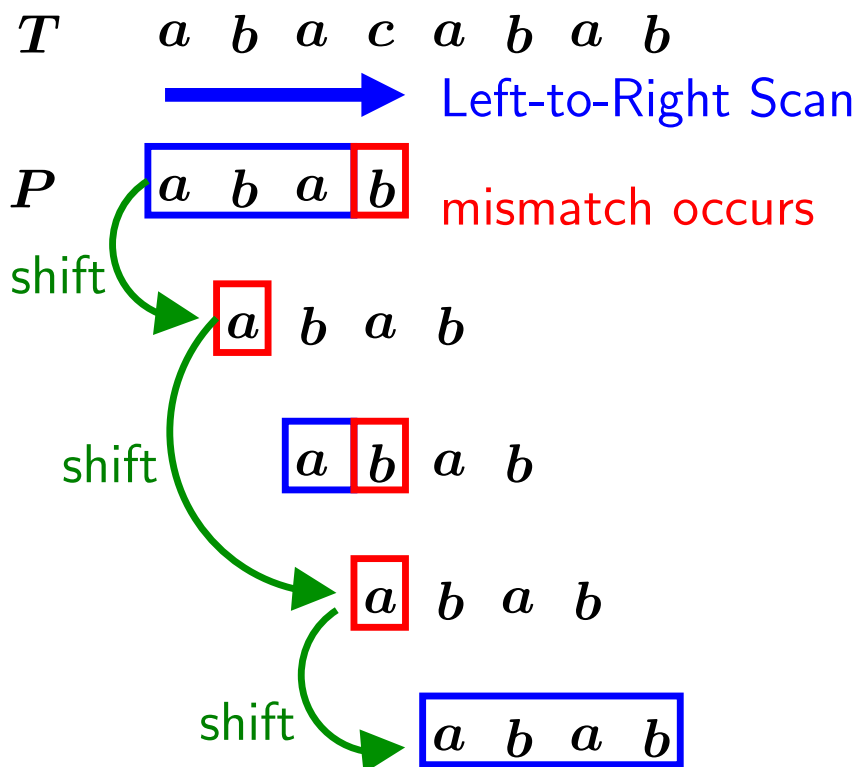   - Time complexity: $\mathcal{O}(m+n)$
2. Boyer-Moore algorithm
   - Proposed by Boyer and Moore at 1977
   - Time complexity: $\mathcal{O}(m+n)$

# KMP algorithm

- Left-to-right scan
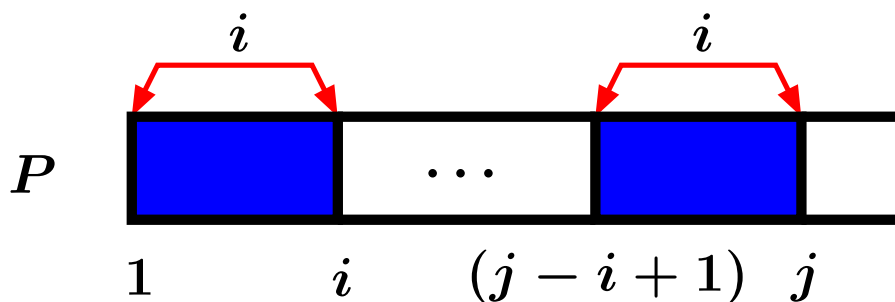- Failure function shift rule

# KMP algorithm

$$T \qquad a \quad b \quad a \quad c \quad a \quad b \quad a \quad b$$

Left-to-Right Scan

$$P \qquad \boxed{a \quad b \quad a} \; \boxed{b}$$ mismatch occurs

shift

$$\boxed{a} \quad b \quad a \quad b$$

shift

$$\boxed{a} \; \boxed{b} \quad a \quad b$$

$$\boxed{a} \quad b \quad a \quad b$$

shift

$$\boxed{a \quad b \quad a \quad b}$$

# Failure function

- If $P = p_1 p_2 \cdots p_n$, then its **failure function** $f$ is defined as follows, where $1 \leq j \leq n$.

$$f(j) = \begin{cases} \text{largest } i < j \text{ such that} & \text{if such an} \\ p_1 \cdots p_i = p_{j-i+1} \cdots p_j, & i \geq 1 \text{ exists}, \\ 0, & \text{otherwise}. \end{cases}$$



$$P$$

$$1 \qquad i \qquad (j-i+1) \quad j$$

# Examples of failure function

- Example 1:

$$P \quad a \quad b \quad a \quad b \quad a \quad b \quad c$$
$$f(j) \quad 0 \quad 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 0$$

- Example 2:

$$P \quad a \quad b \quad a \quad c \quad a \quad b \quad a \quad b$$
$$f(j) \quad 0 \quad 0 \quad 1 \quad 0 \quad 1 \quad 2 \quad 3 \quad 2$$

# KMP algorithm

$/^* \; T = T_1 T_2 \cdots T_m$ and $P = p_1 p_2 \cdots p_n \; ^*/$

```
1  i = 1; j = 1;
2  while i ≤ m and j ≤ n do
3      if Tᵢ = pⱼ then
4          i = i + 1; j = j + 1;
5      else if j = 1 then i = i + 1; j = 1;
6          else j = f(j − 1) + 1;
7      end if
8  end while
9  if j = n + 1 then "a match" else "no match".
```

# KMP algorithm

```
        1  2  3  4  5  6  7  8
T       a  b  a  c  a  b  a  b
f(j)    0  0  1  2
```

$P$   $\boxed{a \ b \ a \ \boxed{b}}$     $(1,1) \rightarrow (4,4), j = f(3) + 1 = 2$

$\boxed{a \ b \ a \ b}$     No this step

$\boxed{a \ \boxed{b}} \ a \ b$     $(4,2) \rightarrow (4,2), j = f(1) + 1 = 1$

$\boxed{a} \ b \ a \ b$     $(4,1) \rightarrow (4,1), i = 4 + 1 = 5, j = 1$

$\boxed{a \ b \ a \ b}$     $(5,1) \rightarrow (8,4), i = 8 + 1, j = 4 + 1$

---

# The correctness of KMP algorithm



Then $p_1 \cdots p_k = T_{j-k} \cdots T_{j-1} = p_{j-k} \cdots p_{j-1}$

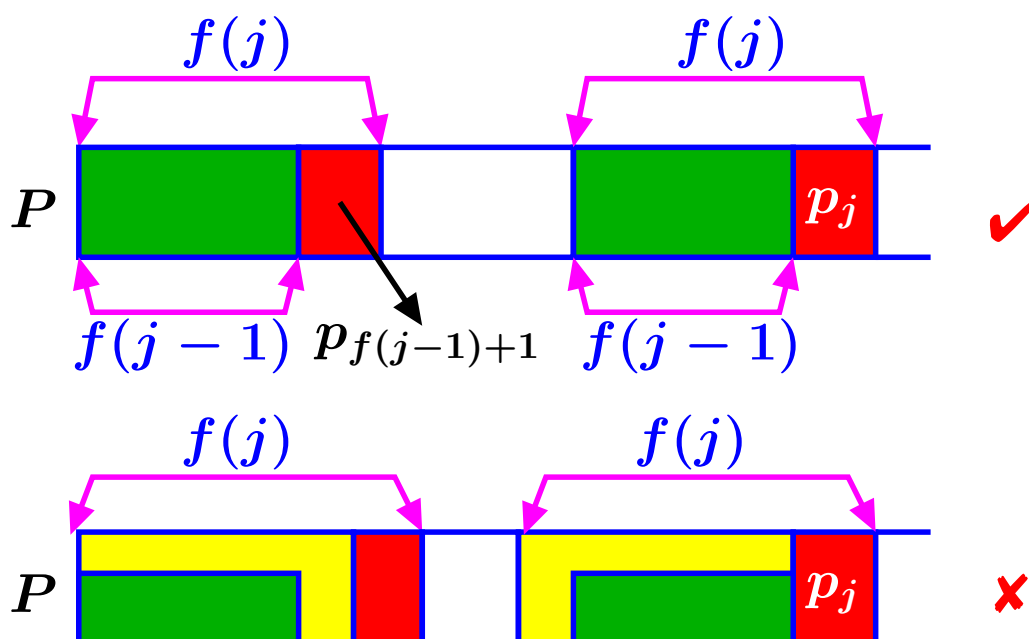Hence $f(j-1) = k > f(j-1)$, a contradiction.

# How to compute failure function?

- Let $f^1(j) = f(j)$ and $f^k(j) = f(f^{k-1}(j))$.
- Let $P = p_1 p_2 \cdots p_n$. Then, $f(1) = 0$ and for $j \geq 2$,

$$f(j) = \begin{cases} f^k(j-1) + 1 & \text{if there exists } k \text{ which is the least integer such that } p_{f^k(j-1)+1} = p_j, \\ \\ 0 & \text{otherwise.} \end{cases}$$
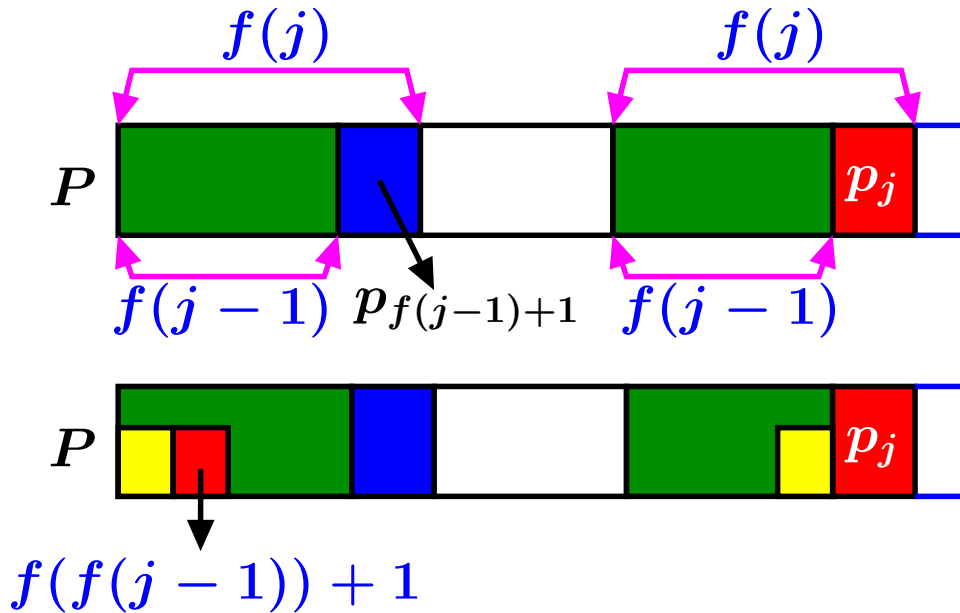
# How to compute failure function?

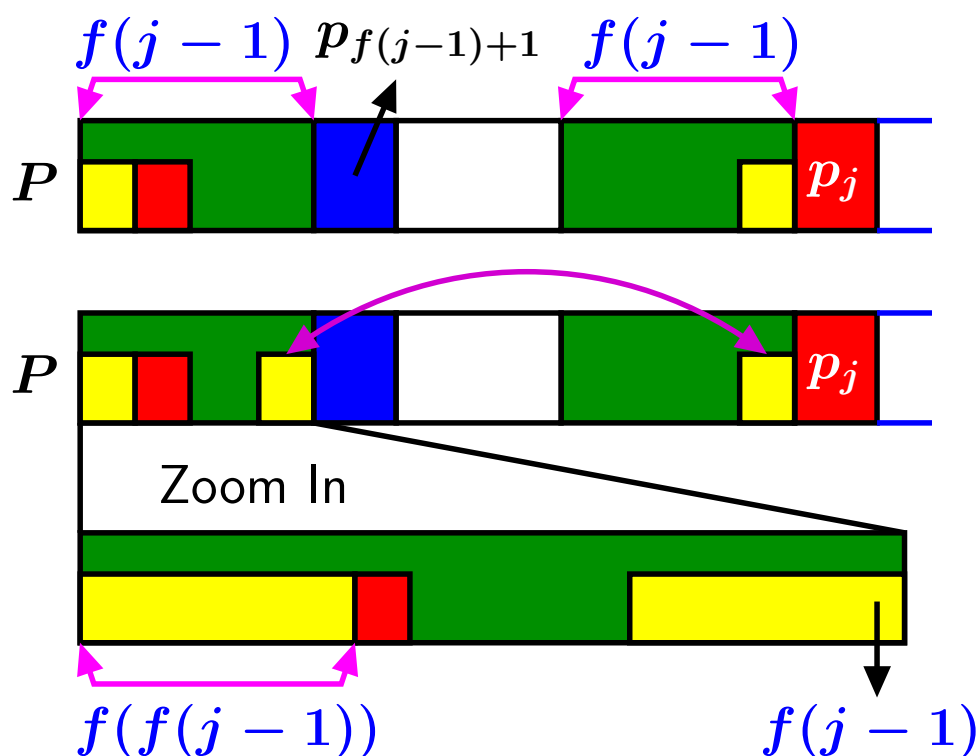- If $p_{f(j-1)+1} = p_j$, then $f(j) = f(j-1) + 1$.

# How to compute failure function?

- If $p_{f(j-1)+1} \neq p_j$ and $p_{f(f(j-1))+1} = p_j$, then $f(j) = f(f(j-1)) + 1$.
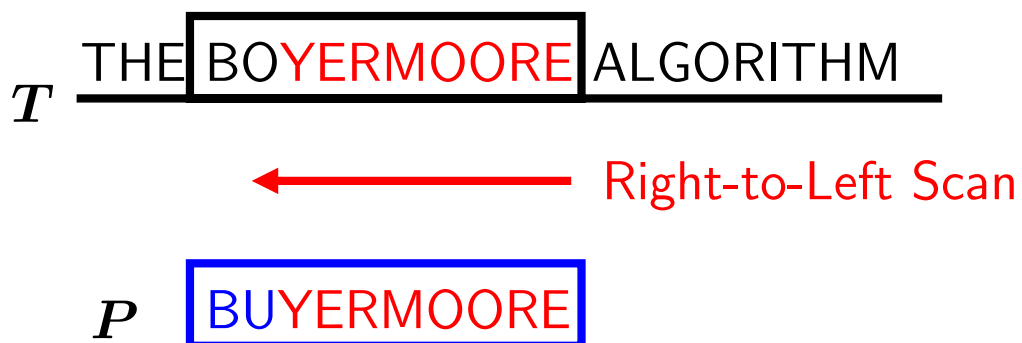
# How to compute failure function?

# Boyer-Moore algorithm

- Right-to-left scan
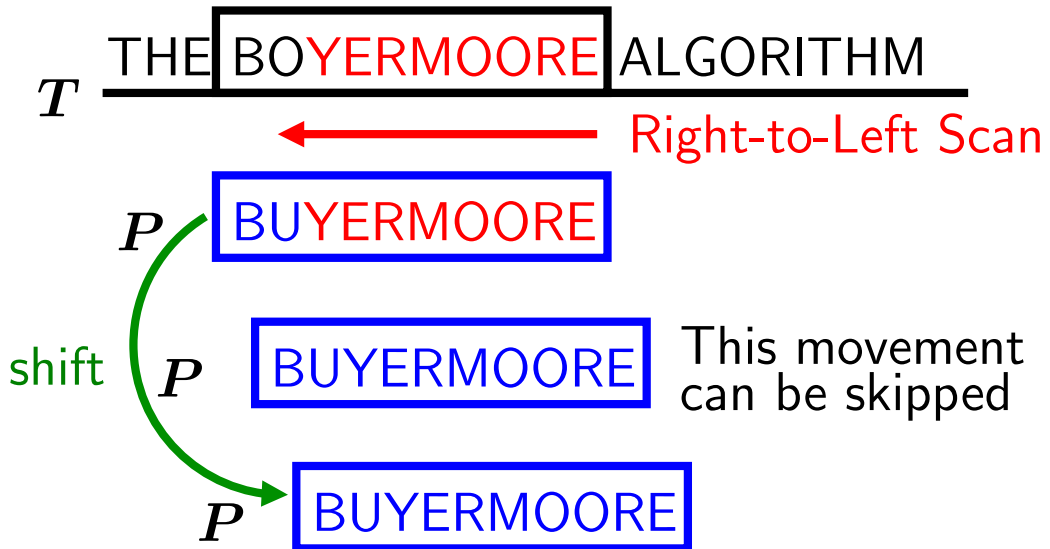- Bad character shift rule
- Good suffix shift rule

# Right-to-left scan

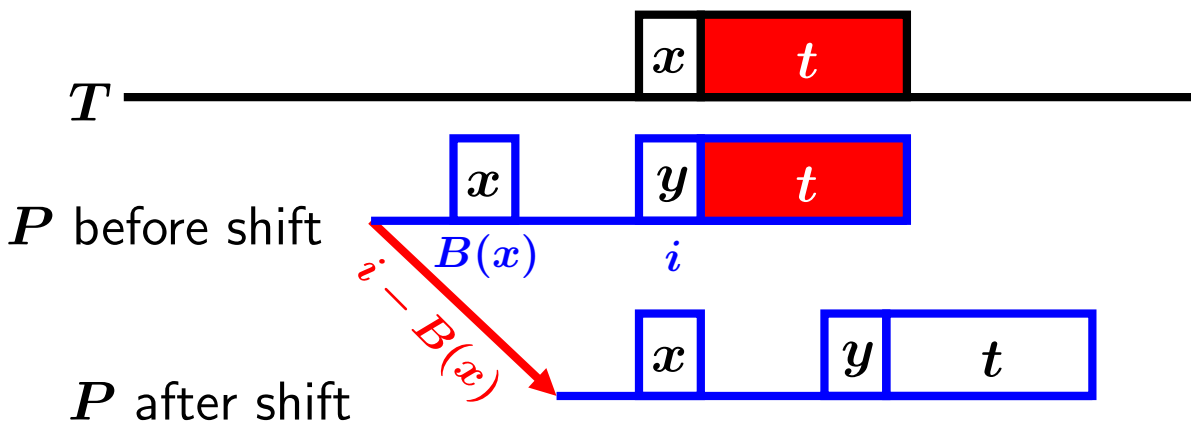- Check whether $P$ occurs in $T$ at some position in the right-to-left scaning manner

$T$    THE BO**YERMOORE** ALGORITHM

⟵ Right-to-Left Scan

$P$    BUYERMOORE

# Bad character shift rule

- What happened if the initial mismatch occurs?

$$T \quad \text{THE BOYERMOORE ALGORITHM}$$

← Right-to-Left Scan

$P$ BUYERMOORE

shift $P$ BUYERMOORE   This movement can be skipped
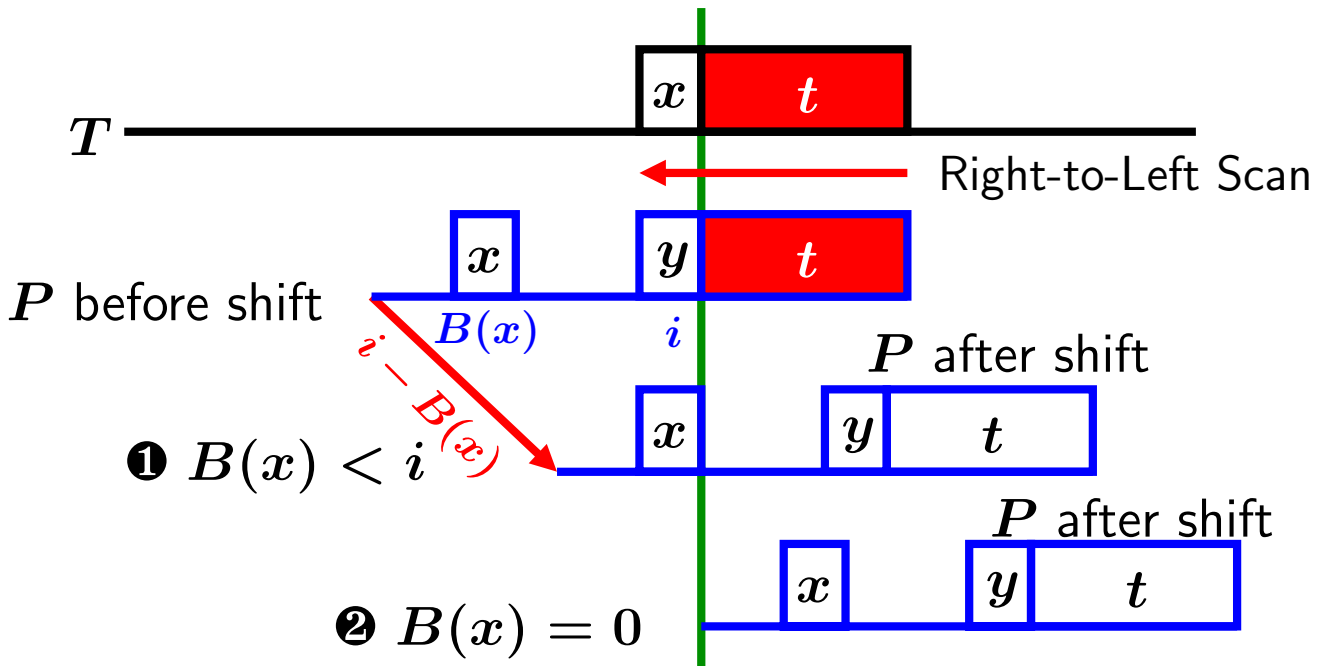
$P$ BUYERMOORE

---

# Bad character shift rule

- Bad character rule: If $x \neq y$, then $P$ is shifted right by $\max\{1, i - B(x)\}$ places.

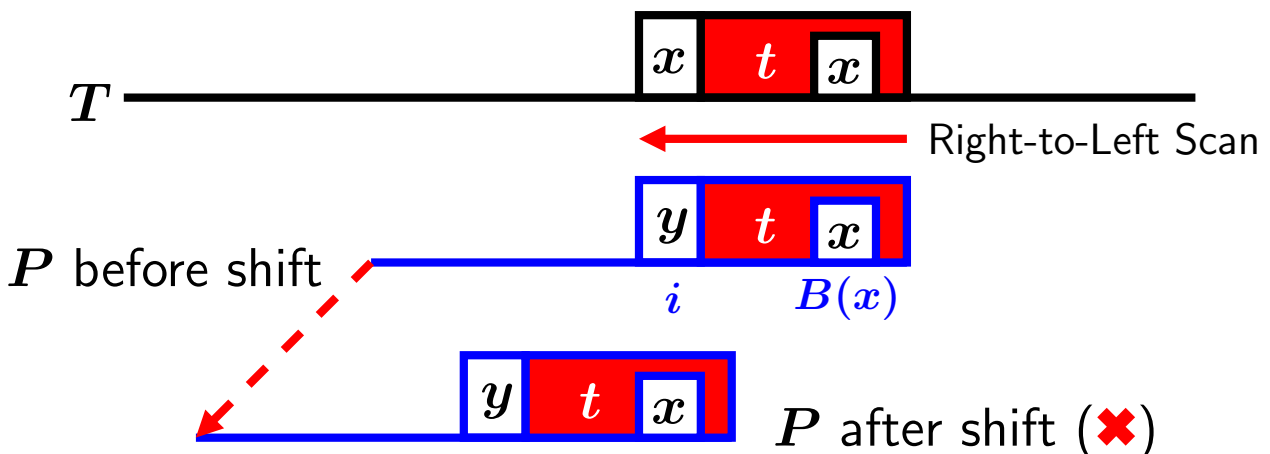- $B(x)$: the position of right-most occurrence of $x$ in $P$ ($B(x) = 0$ if $x$ does not occur in $P$)

$$T \quad \boxed{x}\ \boxed{t}$$

$P$ before shift $\quad \boxed{x}\quad \boxed{y}\ \boxed{t}$

$B(x) \qquad i$

$i - B(x)$

$P$ after shift $\quad \boxed{x}\quad \boxed{y}\ \boxed{t}$

# Bad character shift rule



$T$

$x$ | $t$

Right-to-Left Scan

$P$ before shift

$x$ | $y$ | $t$

$B(x)$ | $i$

$i - B(x)$

$P$ after shift

❶ $B(x) < i$

$x$ | $y$ | $t$

$P$ after shift

❷ $B(x) = 0$

$x$ | $y$ | $t$

# Bad character shift rule

- If $B(x) > i$, then bad character rule has no effect.



$T$

$x$ | $t$ | $x$

Right-to-Left Scan

$P$ before shift

$y$ | $t$ | $x$

$i$ | $B(x)$

$y$ | $t$ | $x$   $P$ after shift (✖)

# Bad character shift rule

- How to compute $B(x)$ for all $x$ in $P$? ($n = |P|$)

  **for** $i = 1$ to $n$ **do**
      $B(x) = 0;$
  **end for**
  **for** $i = n$ to $1$ **do**
      **if** $B(P[i]) = 0$ **then** $B(P[i]) = i;$
  **end for**

- **Extended** bad character shift rule:
  For each position $i$ in $P$ and for each $x$ in $\Sigma$, find the position of the closest occurrence of $x$ in $P$ to the left of $i$.
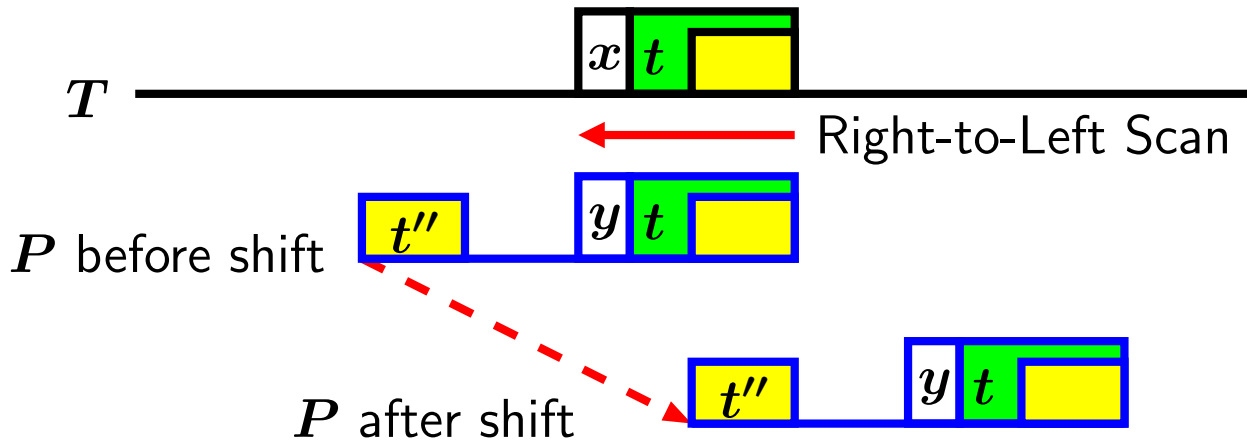
# Good suffix rule: Case ❶

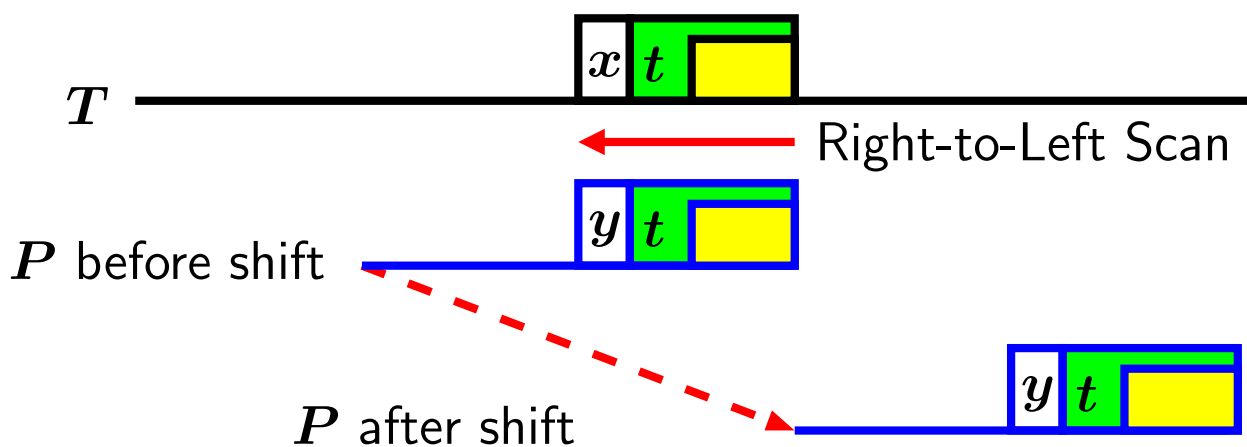- If $x \neq y$, then find the right-most copy $t'$ of $t$ in $P$ such that $t'$ is not a suffix of $P$ and $z \neq y$

# Good suffix rule: Case ❷

- If $t'$ does not exist, then find the largest prefix $t''$ of $P$ such that it is equal to a suffix of $t$

$T$ ─── $x$ $t$ [ ]

Right-to-Left Scan ←

$P$ before shift: $t''$ $y$ $t$ [ ]

$P$ after shift: $t''$ $y$ $t$ [ ]

---

# Good suffix rule: Case ❸

- If $t'$ and $t''$ do not exist, then

$T$ ─── $x$ $t$ [ ]

Right-to-Left Scan ←

$P$ before shift: $y$ $t$ [ ]

$P$ after shift: $y$ $t$ [ ]

# Good suffix rule

$T$

$x$ | $t$

Right-to-Left Scan

$z$ | $t'$    $y$ | $t$    $P$ before shift

$P$ after shift for case ❶

$z$ | $t'$    $y$ | $t$

$P$ after shift for case ❷

$t''$    $y$ | $t$

$P$ after shift for case ❸

$y$ | $t$