# Ch2. Combinatorial Pattern Matching Algorithms
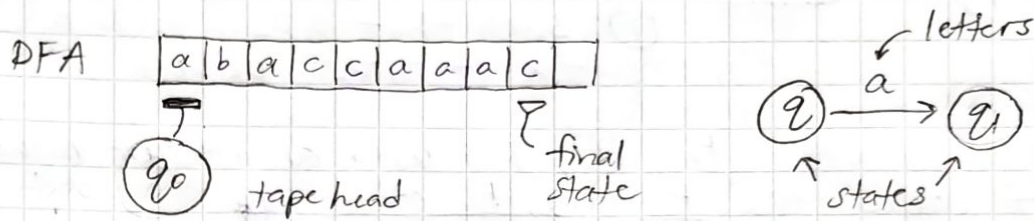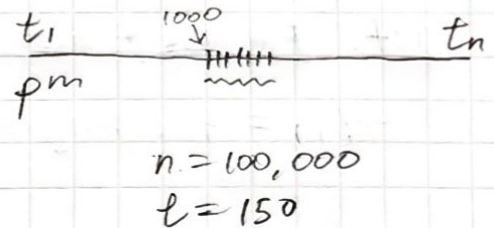
2.1. DFA, NFA, Regular Expressions
2.2. Knuth-Morris-Pratt Algorithm (KMP)

DFA

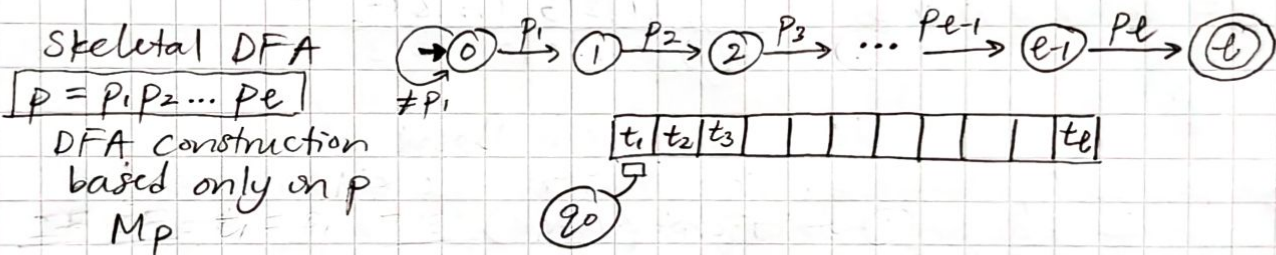| $a$ | $b$ | $a$ | $c$ | $c$ | $a$ | $a$ | $a$ | $c$ | | |

$q_0$    tape head     final state

letters

$q$ $\xrightarrow{a}$ $q_1$

states

Pb.    **Input**    text $t = t_1 t_2 \ldots t_n$
         pattern $p = p_1 p_2 \ldots p_\ell$
         $n \gg \ell$
         $t_i, p_j \in A$

$t_1$      1000        $t_n$
$p^m$

$n = 100,000$
$\ell = 150$

   **Output**    The position in $t$ of
         the first occurrence of $p$

## 2.2. KMP Algorithm

Skeletal DFA

$\boxed{p = p_1 p_2 \ldots p_\ell}$

$\xrightarrow{} \boxed{0} \xrightarrow{p_1} \boxed{1} \xrightarrow{p_2} \boxed{2} \xrightarrow{p_3} \ldots \xrightarrow{p_{\ell-1}} \boxed{\ell-1} \xrightarrow{p_\ell} \boxed{\ell}$
     $\neq p_1$

DFA construction
based only on $p$
$M_p$

| $t_1$ | $t_2$ | $t_3$ | | | | | | | $t_\ell$ |

$q_0$

   **IF** $t_1 = p_1$    **THEN**    $M_p$ enters state ①
                            The head moves to $t_2$

   **IF** $t_1 \neq p_1$    **THEN**    $M_p$ stays in state ⓪
                            the head moves to $t_2$

an inductive argument
"base case"

Suppose after having "read" $t_1 t_2 \ldots t_k$ we have $M_p$ in
state ①. This implies that the last $j$ letters of
$t_1 t_2 \ldots t_k$ are $p_1 p_2 \ldots p_j$.

     $p_1 p_2 \ldots p_j$      The suffix of the prefix of length $k$ of $t$
                         is the prefix of length $j$ of $p$

Induction hypothesis

"inductive step"

IF $t_{k+1} = P_{j+1}$ THEN $M_p$ enters state $(j+1)$ and advances head to $t_{k+2}$

IF $t_{k+1} \neq P_{j+1}$ THEN $M_p$ enters the <u>highest number state</u> $(i)$ such that $P_1 P_2 \ldots P_i$ is a suffix of $t_1 t_2 \ldots t_k t_{k+1}$

To help with discovering this $(i)$, $M_p$ uses an integer valued function $f$ called the <u>Failure Function</u> $(*)$

Def. <u>Failure Function</u> $f: \{1, 2, 3, \ldots, \ell\}$

$f(j) =$ the largest $\boxed{s < j}$ such that $P_1 P_2 \ldots P_s$ is a suffix of $P_1 P_2 \ldots P_j$

$$P_1 P_2 \ldots P_s = P_{j-s+1} \ldots P_{j-1} P_j$$

otherwise $f(j) = 0$

$p = a a b b a a b$

| $j$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| $f(j)$ | 0 | 1 | 0 | 0 | 1 | 2 | 3 |

$(*)$ We will present an algorithm for the construction of the Failure Function. (later)

To see how the failure function is used by $M_p$, let us define $f^{(m)}(j)$ as follows:

i) $f^{(1)}(j) = f(j)$
ii) $f^{(m)}(j) = f(f^{(m-1)}(j))$    for $m > 1$
   i.e. $f^{(m)}(j) = \underbrace{f(f(f \ldots f(j)))}_{m}$

$f(6) = 2, \quad f(2) = 1 \implies f^{(2)}(6) = 1$

Suppose again $M_p$ is in state $(j)$ having read $t_1 t_2 \ldots t_k$ and $t_{k+1} \neq P_{j+1}$. At this point $M_p$ applies the failure function <u>repeatedly</u> to $(j)$ until it finds the smallest value of $m$ for which either

<u>Case 1</u> $f^{(m)}(j) = u$ and $t_{k+1} = P_{u+1}$    or

<u>case 2</u> $f^{(m)}(j) = 0$ and $t_{k+1} \neq P_1$

That is, $M_p$ backs up through states $f^{(1)}(j)$, $f(f^{(1)}(j)) = f^{(2)}(j)$, ... until __Case 1__ or __Case 2__ shows up for $f^{(m)}(j)$ but not for $f^{(m-1)}(j)$.

In __Case 1__, $M_p$ enters state (u+1)
In __Case 2__, $M_p$ enters state ⓪
In either case, the head is advanced to $t_{k+1}$

In __Case 1__, it is easy to verify that if $p_1 p_2 ... p_j$ was the longest prefix of $p$ that is a suffix of $t_1 t_2 ... t_k$, then $p_1 p_2 ... p_{f^{(u)}(j)+1}$ is the longest prefix of $p$ that is a suffix of $t_1 t_2 ... t_k t_{k+1}$

In __Case 2__, no prefix of $p$ is a suffix of $t_1 t_2 ... t_k t_{k+1}$

$M_p$ then reads $t_{k+2}$. $M_p$ continues operating in this fashion __either__ until it enters the final state ⓔ, in which case we know that the last input symbol of gives us a complete instance of pattern $p = p_1 p_2 ... p_e$, __or__ until $M_p$ has read the last symbol of $t$ without entering final state ⓔ.

$p = a\ a\ b\ b\ a\ a\ b$       $A = \{a, b\}$
$t = a\ b\ a\ a\ b\ a\ a\ b\ b\ a\ a\ b$