

Real-time string matching in sublinear space

Leszek Gąsieniec, Roman Kolpakov
*Department of Computer Science,
University of Liverpool, Liverpool, UK*

String matching problem

$T = t_1 \dots t_n$ — text

$P = p_1 \dots p_m$ — pattern ($m \leq n$)

The problem is to find all occurrences of P in T

Ex:

$P = a b a a$

$T = a b a b a a b a a b b a b a a b$



occurences of P in T

Next function

$$P = p_1 \dots p_m$$

$$\text{next} : \{1, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$$

$\text{next}(i)$ is the largest positive k s.t.

$p_1 \dots p_{k-1} = p_{i-k+1} \dots p_{i-1}$ and $p_i \neq p_k$
if such k exists ($\text{next}(i) = 0$ otherwise).

$$P: \begin{array}{cccc} p_1 & p_{k-1}p_k & p_{i-k+1} & p_{i-1}p_i \\ \hline \hline \text{next}(i) = k & & & p_k \neq p_i \end{array}$$

Ex: $P = abaaababababc$

$$\text{next}(1) = 0; \text{next}(2) = 1; \text{next}(3) = 0;$$

$$\text{next}(4) = 2; \text{next}(5) = 2; \text{next}(6) = 1;$$

$$\text{next}(7) = 0; \text{next}(8) = 4; \text{next}(9) = 0;$$

$$\text{next}(10) = 4; \text{next}(11) = 0;$$

$$\text{next}(12) = 4; \text{next}(13) = 3$$

next can be computed in $O(m)$ time

Knuth-Morris-Pratt algorithm

KMP algorithm:

compute **next** for $p_1 \dots p_m$

$i := 1$; {pointer to pattern symbol p_i }

for $j := 1$ **to** n {pointer to text symbol t_j }

while $i > 0$ **and** $p_i \neq t_j$

$i := \mathbf{next}(i)$;

if $i = m$

 report an occurrence of P ending in t_j ;

$i := \mathbf{next}(m + 1)$;

else $i := i + 1$;

$T : \underline{t_1 \quad t_{j-i+1} \quad t_j}$

$P : \underline{p_1 \quad p_k \quad p_i} \quad p_i \neq t_j$

$k = \mathbf{next}(i) \quad P : \underline{p_1 \quad p_k}$

Real-time string matching

Def: *On-line* string matching — input text symbols are assessed consecutively and any occurrence of pattern is reported during assessing the last symbol of the occurrence

Def: *Real-time* string matching — on-line string matching such that the time for assessment of an input text symbol is bounded by a constant

$f_0 = b, f_1 = a, f_n = f_{n-1}f_{n-2}$ for $n \geq 2$

$|f_n| = F_n$ — n -th Fibonacci number

f_n and $f_{n-2}f_{n-1}$ coincide in the first $F_n - 2$ letters and don't coincide in the $(F_n - 1)$ -th letter

Ex: $f_5 = \underline{abaababa}$

$f_3f_4 = \underline{abaabaab}$

if $P = f_n$ then **next** $(F_n - 1) = F_{n-1} - 1$

⇓

the **while**-loop in KMP algorithm can be executed at least n times for the same text symbol t_j

⇓

KMP algorithm is not a real-time algorithm

Real-time KMP modification

compute **next** for $p_1 \dots p_m$ \$

$i := 1; j' = 1;$

for $j := 1$ **to** n

 KMP_STEP();

if $j' \leq j$ **then** KMP_STEP();

KMP_STEP():

if $p_i = t_{j'}$

if $i = m$ **then**

 report an occurrence of P ending in $t_{j'}$;

$i := \text{next}(m + 1);$

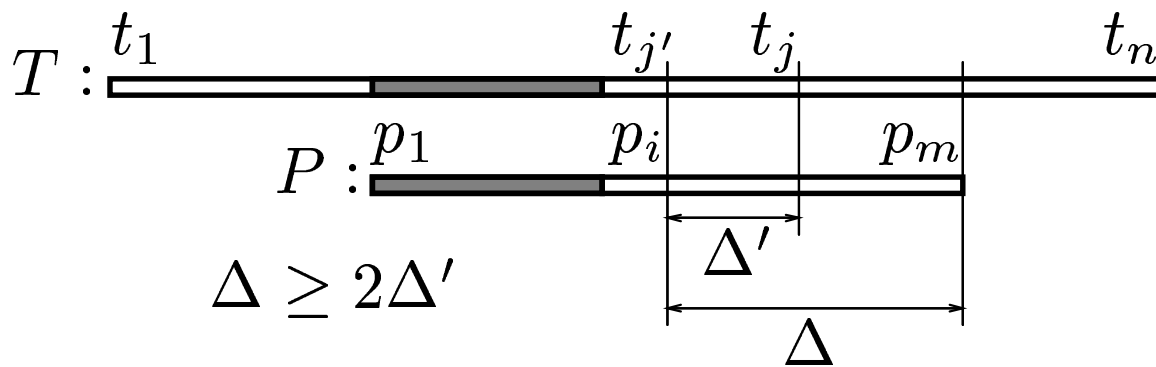
else $i := i + 1;$

$j' := j' + 1;$

else $\{p_i \neq t_{j'}\}$

if $\text{next}(i) > 0$ **then** $i := \text{next}(i);$

else $j' := j' + 1; i := 1;$



a) the algorithm finds any occurrence $T[u..v]$ of P during the processing the end symbol t_v ;

b) $\Delta' \leq m/2$.

Galil 1981:

The algorithm is a real-time algorithm which solves the string matching problem in $O(m + n)$ time and $O(m)$ space

Space efficient string matching

Can string matching problem be solved in $o(m)$ space?

Galil, Seiferas 1983:

Crochemore, Perrin 1991:

Gąsieniec, Plandowski, Rytter 1995:

String matching problem can be solved in $O(m + n)$ time and $O(1)$ space

Gąsieniec, Potapov 2003:

On-line algorithm which solves the string matching problem in $O(m + n)$ time and $O(1)$ space

All the algorithms are not real-time

Our problem

Can real-time string matching problem
be solved in $o(m)$ space?

Partial next function

for any real x

$$\mathbf{next}_x(i) = \begin{cases} \mathbf{next}(i) & \text{if } \mathbf{next}(i) > x, \\ 0 & \text{otherwise.} \end{cases}$$

Ex: $P = abaaababababc$

$\mathbf{next}_1(1) = 0; \quad \mathbf{next}_1(2) = 0;$

$\mathbf{next}_1(3) = 0; \quad \mathbf{next}_1(4) = 2;$

$\mathbf{next}_1(5) = 2; \quad \mathbf{next}_1(6) = 0;$

$\mathbf{next}_1(7) = 0; \quad \mathbf{next}_1(8) = 4;$

$\mathbf{next}_1(9) = 0; \quad \mathbf{next}_1(10) = 4;$

$\mathbf{next}_1(11) = 0; \quad \mathbf{next}_1(12) = 4;$

$\mathbf{next}_1(13) = 3$

list[**next**₁]:

$$(4, 2)(5, 2)(8, 4)(10, 4)(12, 4)(13, 3) =$$
$$\underbrace{(4, 2)(5, 2)}_{\alpha_1} \underbrace{(8, 4)(10, 4)(12, 4)}_{\alpha_2} \underbrace{(13, 3)}_{\alpha_3}$$

list*[**next**₁] = $\alpha_1\alpha_2\alpha_3$ — list of ranges

$$\mathbf{value}(\alpha_2) = 4, \quad \mathbf{begin}(\alpha_2) = 8,$$

$$\mathbf{end}(\alpha_2) = 12, \quad \mathbf{step}(\alpha_2) = 2$$

⇓

each range requires $O(1)$ space.

Lemma:

list*[**next** _{x}] contains $O(m/x)$ ranges and requires $O(m/x)$ space.

Let $\tau \geq 2$

$m_i = \lfloor m^{i/\tau} \rfloor$, $P_i = P[1..m_i]$, $i = 1, \dots, \tau$

h_0 — **next** function for P_1 \$

h_i — **next** _{$m^{i/\tau}$} function for P_{i+1} \$,

$i = 1, \dots, \tau - 1$

Ex: $P = abaaabababab$, $\tau = 2$

$m = |P| = 12$, $m_1 = 3$, $P_1 = aba$ \$

$h_0(1) = 0$, $h_0(2) = 1$

$h_0(3) = 0$, $h_0(4) = 2$

list^{*} $[h_1] = \alpha_2$

h_0 , **list**^{*} $[h_1], \dots, \mathbf{list}^*[h_{\tau-1}]$ require at total of $O(m^{1/\tau} \cdot \tau)$ space.

KMP_STEP computes for each j
largest k , s.t., $p_1 \dots p_k = t_{j-k+1} \dots t_j$

maxpref₀(j) — largest $k \leq m^{1/\tau}$, s.t.,
 $p_1 \dots p_k = t_{j-k+1} \dots t_j$

(**maxpref**₀(j) = 0 if there is no such k)

maxpref _{i} (j), $i = 1, \dots, \tau - 1$, — largest
 k , s.t., $p_1 \dots p_k = t_{j-k+1} \dots t_j$ and
 $m^{i/\tau} < k \leq m^{(i+1)/\tau}$ (**maxpref** _{i} (j) = 0 if
there is no such k)

COMP_MPREF _{i} , $i = 0, 1, \dots, \tau - 1$, —
modification of KMP_STEP for
computing of **maxpref** _{i} , using **list**^{*}[h_i]
(h_0 if $i = 0$)

Space efficient algorithm

compute $h_0, \mathbf{list}^*[h_1], \dots, \mathbf{list}^*[h_{\tau-1}]$

$j_0 := j_1 := \dots j_{\tau-1} := 1;$

$\{j_i - \text{current position for which } \mathbf{maxpref}_i \text{ is computed}\}$

$l_0 := 1; l_1, \dots, l_{\tau-1}$ are undefined;

$\{l_i - \text{position of pattern symbol compared with } t_{j_i}\}$

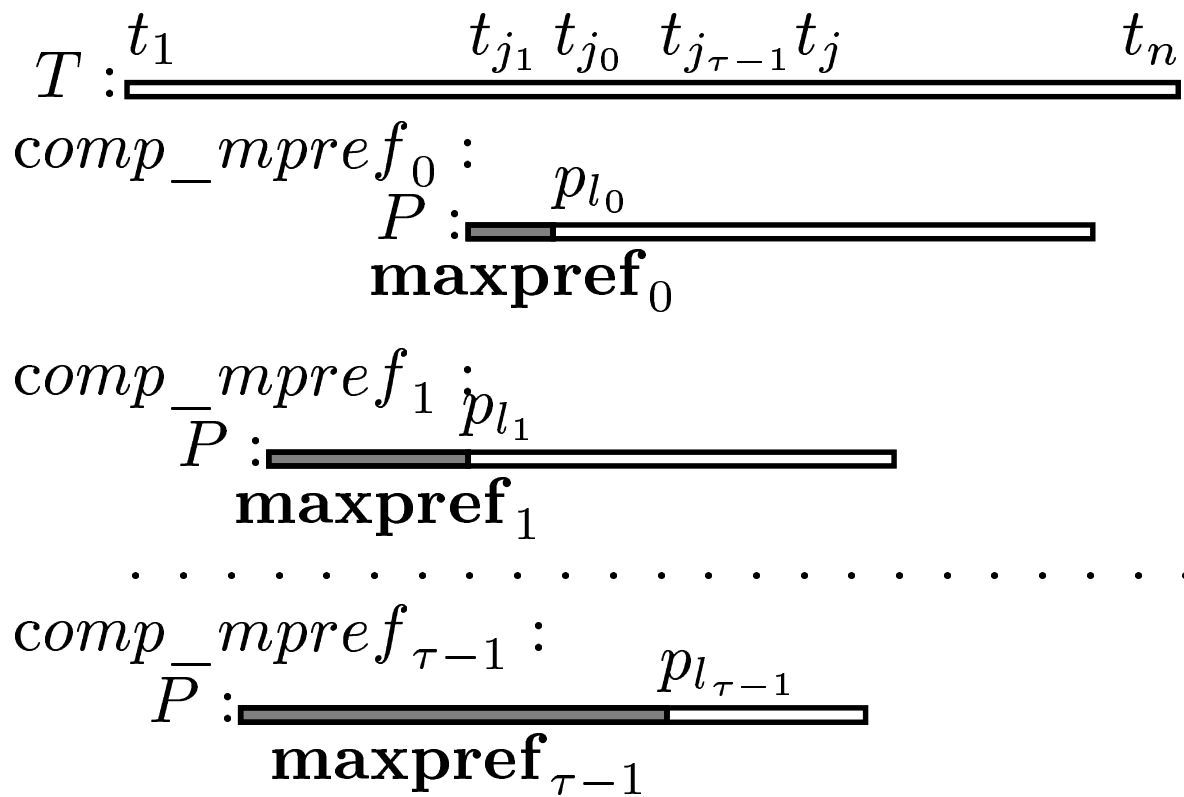
for $j := 1$ **to** n

 read the next symbol t_j ;

for $i := \tau - 1$ **down to** 0

 COMP_MPREF _{i} (\cdot);

if $j_i \leq j$ **then** COMP_MPREF _{i} (\cdot);



$$j - j_i \leq (|\mathbf{list}^*[h_i]| + 1)/2$$

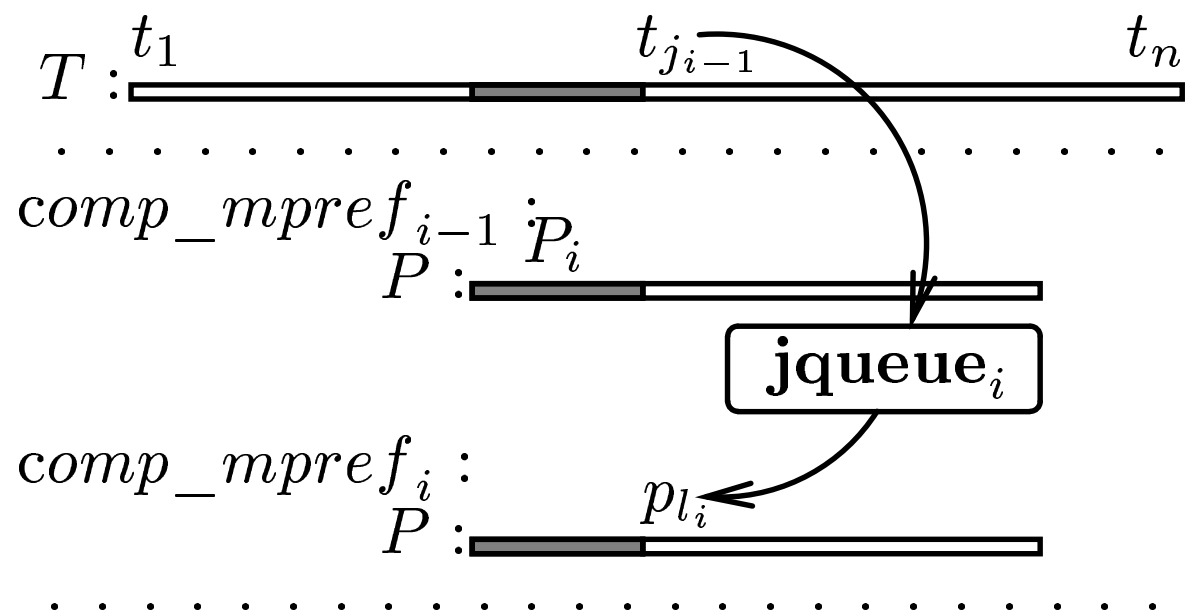
$$\Downarrow$$

$$\max_i (j - j_i) = O(m^{1/\tau})$$

COMP_MPREF_{*i*} begins computation of **maxpref_{*i*}** only after

COMP_MPREF_{*i-1*} finds an occurrence of P_i in T .

Positions of occurrences of P_i are communicated to COMP_MPREF_{*i*} by means of FIFO queue **jqueue_{*i*}**:



$$j - j_i = O(m^{1/\tau}) \Rightarrow |\mathbf{jqueue}_i| = O(m^{1/\tau})$$

Preprocessing

$P : \underline{p_1 \quad p_{j_1} \quad p_{j_0} \quad p_{j_{\tau-1}} \quad p_m}$

$P : \underline{p_{l_0}}$
maxpref₀

$P : \underline{p_{l_1}}$
maxpref₁

.....

$P : \underline{p_{l_{\tau-1}}}$
maxpref _{$\tau-1$}

Preprocessing requires $O(m)$ time and
 $O(m^{1/\tau} \cdot \tau)$ space

COMP_MPREF _{$\tau-1$} reports occurrences of $P_\tau = P$ in T

The algorithm solves the string matching problem in $O(m + n)$ time and $O(m^{1/\tau} \cdot \tau)$ space, and time required for the test of each text symbol is bounded by $O(\tau)$.

Theorem: For any constant $\varepsilon > 0$ the string matching problem can be solved in real time and $O(m^\varepsilon)$ additional space.

Conclusions

Another version of algorithm:

$$j - j_i \leq Km^{1/\tau}, K < m^{1/\tau}$$

$$m'_i = \left\lceil m^{i/\tau} + \frac{Km^{2/\tau}}{m^{1/\tau} - K} \right\rceil, P'_i = P[1..m'_i]$$

$$h'_i = \mathbf{next}_{m^{i/\tau}} \text{ for } P'_{i+1}$$

COMP_MPREF_{*i*} computes

maxpref'_{*i*}(*j*) which is largest *k*, s.t.,

$$p_1 \dots p_k = t_{j-k+1} \dots t_j \text{ and}$$

$$m^{i/\tau} < k \leq m'_{i+1} \quad (k \leq m'_1 \text{ for } i = 0)$$

In this version we can avoid the use of

jqueue_{*i*}.

Is this version is more space efficient?

Can string matching problem be solved in real time and $O(\log m)$ additional space?