# KMP - Algorithm

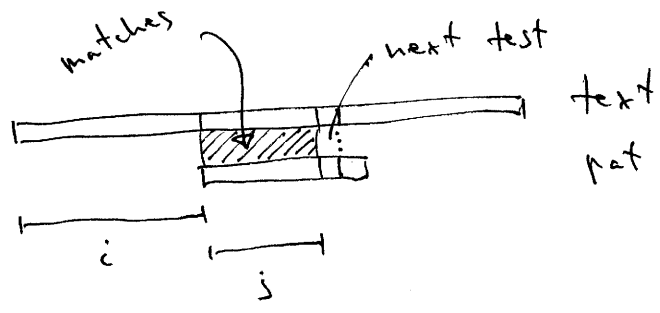$\leftarrow$ ( For reporting all occurences of pat in text )

An Extension of the brute-force alg.

Consider checking pat at shift $i$ in text, where $j$ chars have been matched so far:



Consider the next test $\text{text}[i+j+1] \overset{?}{=} \text{pat}[j+1]$.

We aim at maintaining the following invariants:

1) For current shift $i$ of pattern, the first $j$ chars matches text.

2) All shifts $< i$ of pat have been tested (and reported if pat matched there).

## Case A : Next test is positive (chars matches).

We then do the update
$$j \rightarrow j+1$$

[ Shift $i$ is still a possible occurence, now with one more char known to match ]

## Case B1 : Next test is negative, and $j = 0$

We then do the update
$$i \rightarrow i+1$$

( Shift $i$ is now known not to be an occ., and shift should increase. We increase the smallest possible step, and keep $j=0$. )
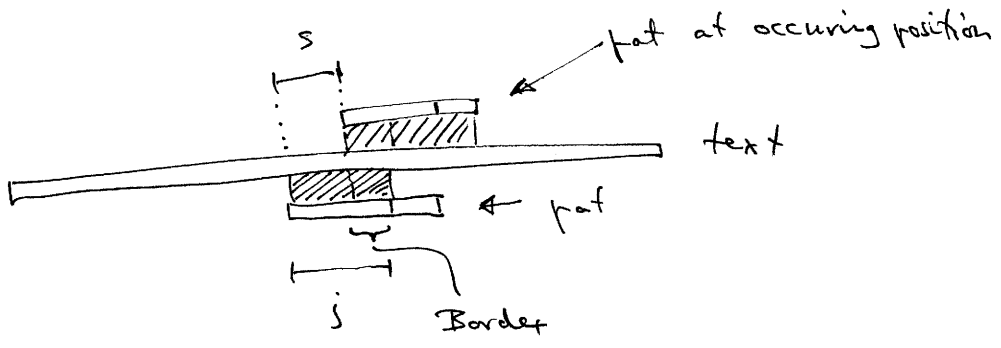
**Case B2** : Next test is negative, and $j \geq 1$

$\left[\begin{array}{l}\text{Shift } i \text{ is now known } \underline{\text{not}} \text{ to be an occ., and} \\ \text{shift should increase. How much can we guaran-} \\ \text{tee it must move ?}\end{array}\right]$

~~Consider the next~~, occ. of pat in text at

Assume there is an

position $i + s$ [ i.e., s further to the right ],

for an $s \leq j$ . [We know $s \geq 1$, so

$1 \leq s \leq j$, which is the reason

why the idea here in B2 does

not hold for $0 = j$, hence

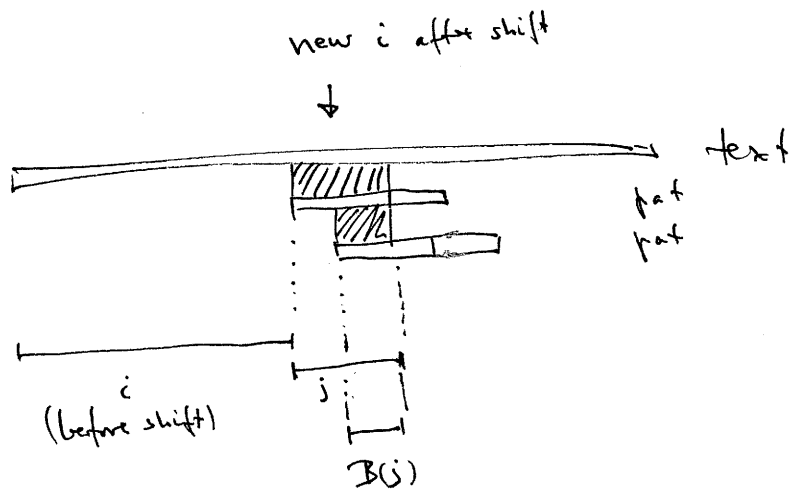the need for the (simple) separate

case B1 }.

That is, we have



Here, we see a border (a suffix which is also a
prefix) of the string pat$[1..j]$ which has length
$j-s$. The border is non-trivial since $s \geq 1$.
$B(j)$ is the length of the longest such border,
pat

So $j-s \leq B(j)$ . [We write $B(j)$ for $B_{pat}(j)$ from now}

$\Updownarrow$

$j - B(j) \leq s$ .

Hence, there can be no occurences the next
$j - B(j)$ shifts, so we can ~~update~~ increase $i$ by this amount.

Doing so leaves $B(j)$ known matches of chars
after the shift [as $B(j)$ corresponds to an _actual_
border] :



We therefore do the updates

$$i \longrightarrow i + (j - B(j))$$
$$j \longrightarrow j - (j - B(j)) = B(j)$$

(and this keeps invariants 1) and 2) ).

---

Note : By the definition $B(0) = -1$, and the observa-
tion $B(k) \geq 0$ for $k \geq 1$   ($B(k)$ is a length)
we can merge cases B1 & B2
into the single update :

$$i \longrightarrow i + (j - B(j))$$
$$j \longrightarrow \max \{ 0, B(j) \}$$

Above, we have not considered checking whether we do lookups in text and pat past their end.

Adding these, the above analysis (and aimed-at invariants) suggests the following algorithm:

$$\underline{KMP\ (text,\ pat)}$$

$$\boxed{\begin{array}{l} n = |text| \\ m = |pat| \end{array}}$$

$i = 0$

$j = 0$

$\underline{while}\ \boxed{i \leq n - m}$

Ensures no lookups past ends

$\underline{if}\ \boxed{j < m \ AND}\ text\,[\,i+j+1\,] == pat\,[\,j+1\,]$

$j = j+1$

$\underline{else}$

Ensures reporting of occurrences (all)

$\boxed{\begin{array}{l} if\ j == m \\ \quad report\ match\ at\ i \end{array}}$

$i = i + (j - B(j))$

$j = max\ \{0,\ B(j)\}$

It is easy to verify that invariants 1) and 2) are maintained, from which correctness follows.

## Complexity :

Let $k = i + j$

In each case, the change $(\Delta)$ to $i, j, k$ are as follows

| Case | $\Delta i$ | $\Delta j$ | $\Delta k$ |
|------|-----|-----|-----|
| A | 0 | 1 | 1 |
| B1 | 1 | 0 | 1 |
| B2 | $\geq 1$ | (?) | 0 |

$$\Delta k = \Delta i + \Delta j$$
$$= (j - B(j)) - (j - B(j))$$
$$= 0$$

$$\Delta i = (j - B(j)) \geq 1 \quad \longleftarrow \text{ as borders are nontrivial}$$
$$(B(j) < j)$$

So for $z = k + i$ we see that :

$$\begin{bmatrix} z = k + i \\ \Downarrow \\ \Delta z = \Delta k + \Delta i \end{bmatrix}$$

i) $z$ $(= k + i = i + j + i)$ is $0$ after initialization i KMP [first two lines]

ii) $\Delta z \geq 1$ for each loop traversed (which falls in $_\wedge$ cases A, B1, or B2).
    exactly one of

Assume the loop is traversed $t$ times.

At entry to t'th loop:

as we $\overset{do}{\underset{\vee}{\text{enter}}}$ loop t [see loop condition]

$$\underset{\underset{\text{already}}{\wedge}}{\#\text{ loops done}} = t-1 \leq z = 2\cdot i + j \leq 2(n-m) + j$$

$$\leq 2(n-m) + m$$

$\uparrow$ j can never exceed m in code. [as $\mathcal{B}(j) < j$].

for last line of code.

$\Updownarrow$

$$\underline{t \leq 2n - m + 1} \quad \left(\begin{array}{l}\text{So alg.}\\ \text{terminates}\end{array}\right)$$

Note: last iteration must use the <u>else</u>-part (as i is not changed in ~~if~~-part and loop condition is based on i).

So <u>#comparisons performed</u> (~~if~~-part) is $\leq 2n - m$.

For the text $\underset{n}{\underbrace{aaa\cdots a}}$ and pattern $\underset{m}{\underbrace{abbb\cdots b}}$ the algorithm <u>will</u> perform $2n-m$ comparisons. Hence, this can be said to be <u>the exact</u> <u>worst case complexity</u> of KMP.

Summing up :

For $|text| = n$ and $|pat| = m$ we have

Brute-Force takes time $\Theta((n-m) \cdot m)$
[ which can be $\Theta(n^2)$ , eg. for
$m = n/2$ ].

KMP takes time $\Theta(n)$ , assuming
the funktion $B(j)$ is available (as
a precomputed table, eg.) . Note
that $B(j)$ depends only on pat,
and can be reused for other
text's.

However, for a single text and pat,
the time for KMP is really
$\Theta(n + preprocessing(m))$ , where

preprocessing(m) designates the time
to find $B(j)$ for all $j$.

There is a straightforward alg. arising
from the definition of $B$ which runs in
time $\Theta(m^2)$ [two for-loops]. Can we do better ?