

# 12.3 Exclusion methods

## Fast expected running time

2016200648

조유나

## Exclusion Method

---

- k-mismatch, k- difference method
  - text T length  $m$ , pattern P length  $n$
  - Worst case  $\theta(km)$
- **Goal**
  - **-faster expected running time than  $\theta(km)$**
  - -for a reasonable range of  $k$
- Baeza-Yates and Perleberg (BYP)
- Chang and Lawler (CL)
- Myers methods

# Partition approach to approximate matching

---

## A. Partition

- partition T (or P) into consecutive regions of a given **length r**

## B. Search Phase

- find length-r intervals of T that could be approximate occurrence of P using **exact matching methods**  
-> *surviving intervals*
- Eliminate as many intervals as possible

## C. Check Phase

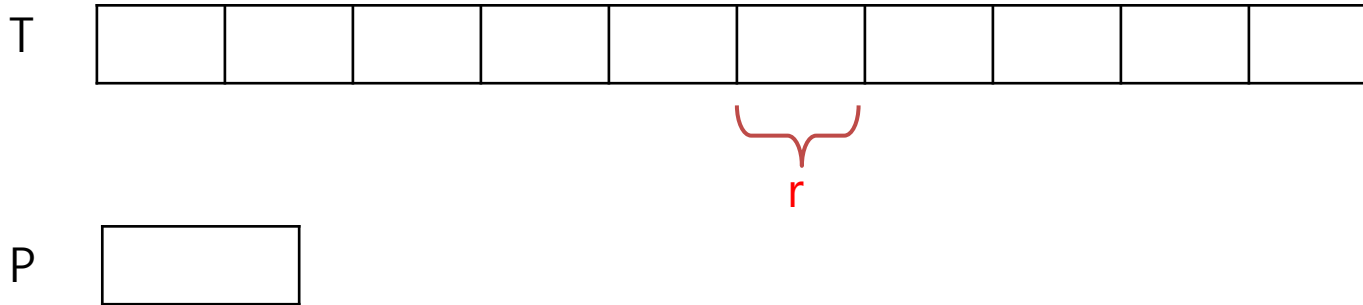
- for each surviving interval R of T, use **approximate matching method** to check approximate occurrence of P in a larger interval R

# Partition approach to approximate matching

---

## A. Partition

- partition T (or P) into consecutive regions of a given **length  $r$**



# Partition approach to approximate matching

---

## B. Search Phase

- find length- $r$  intervals of  $T$  that could be approximate occurrence of  $P$  using **exact matching methods**  
-> *surviving intervals*
- Eliminate as many intervals as possible

Surviving

Non-surviving

T



P

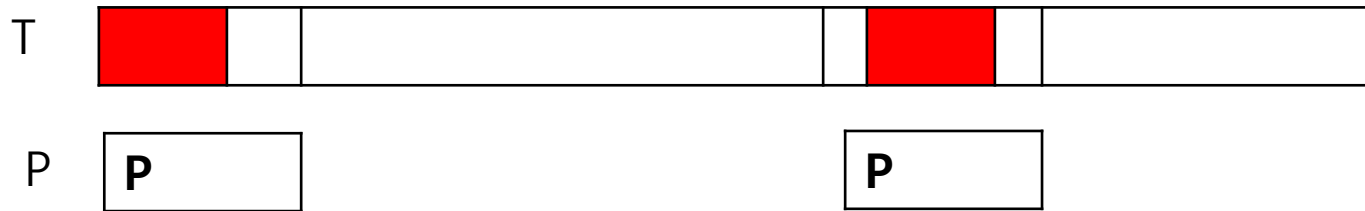


## Partition approach to approximate matching

---

### C. Check Phase

- for each surviving interval R of T, use **approximate matching** method to check approximate occurrence of P in a larger interval R

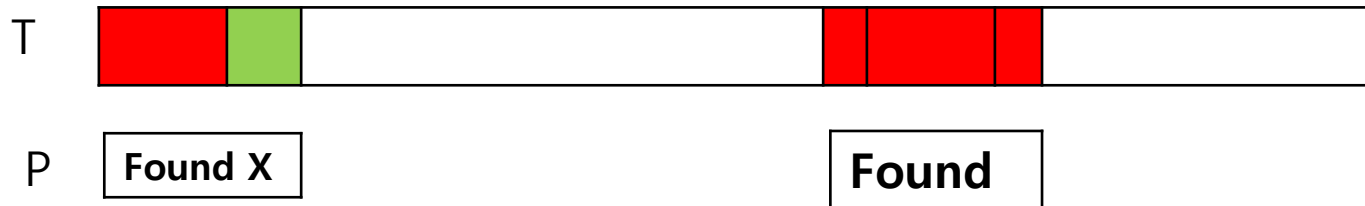


## Partition approach to approximate matching

---

### C. Check Phase

- for each surviving interval R of T, use **approximate matching** method to check approximate occurrence of P in a larger interval R



# Partition approach to approximate matching

---

- **Method differs**
  - choice of  $r$
  - choice of string to partition
  - exact matching method
  - definition of a region in check phase



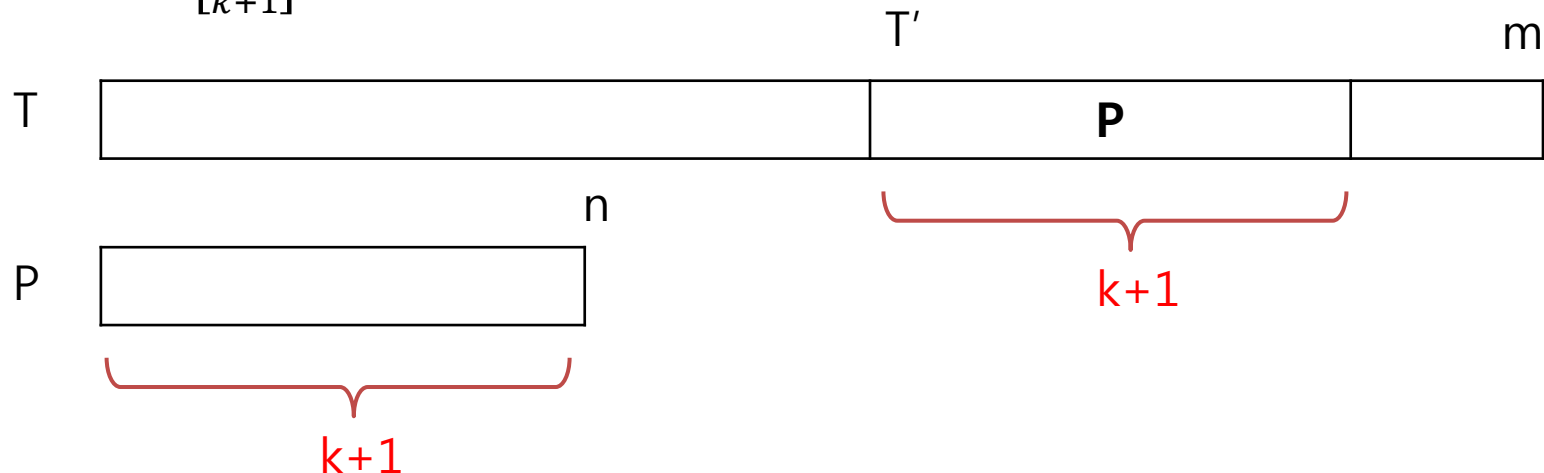
- Baeza-Yates and Perleberg
- $O(m)$  for for  $k < (\frac{n}{\log n})$
- Partition  $r = \left\lfloor \frac{n}{k+1} \right\rfloor$

# BYP method

- **Lemma 12.3.1**

- Suppose  $P$  matches a substring  $T'$  of  $T$  with at most  $k$  differences. Then  $T'$  must contain at least one interval of length  $r$  that exactly matches one of the  $r$ -length regions of the partition of  $P$ .

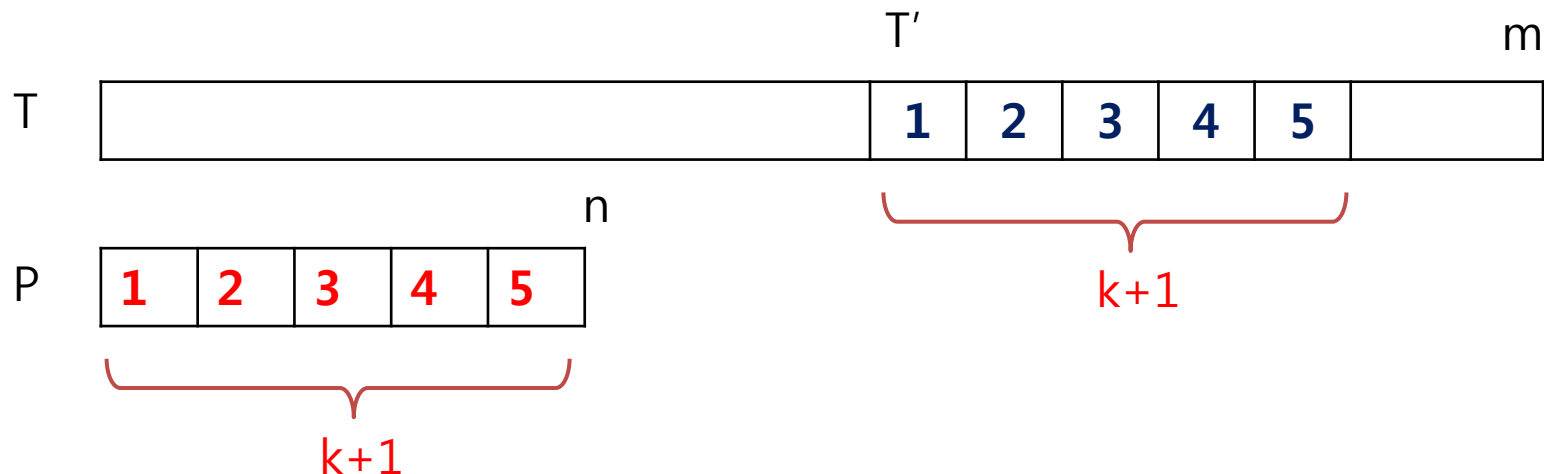
- $r = \left\lfloor \frac{n}{k+1} \right\rfloor$



## BYP method

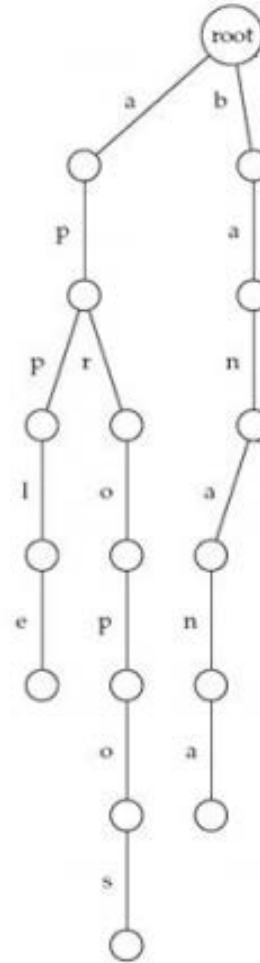
- **Lemma 12.3.1**

- Suppose  $P$  matches a substring  $T'$  of  $T$  with at most  $k$  differences. Then  $T'$  must contain at least one interval of length  $r$  that exactly matches one of the  $r$ -length regions of the partition of  $P$ .
- $k = 4$



# BYP method

- Keyword tree
  - edge : element
- Aho-Corasick algorithm
  - String searching algorithm
  - time complexity  $O(m)$



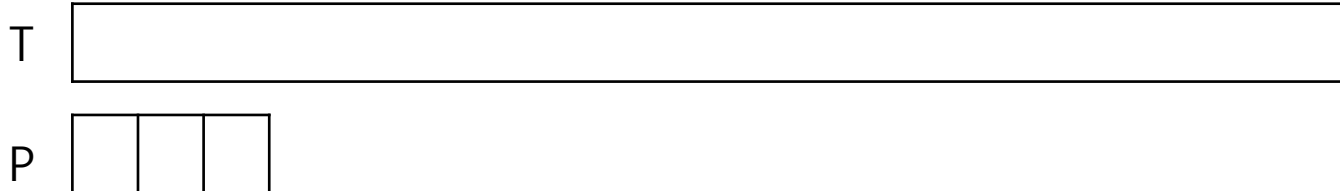
## BYP method

---

- Partition

-Let  $\mathcal{P}$  be the set of  $k+1$  substrings of  $P$  taken from the first  $k+1$  regions of  $P$ 's partition.

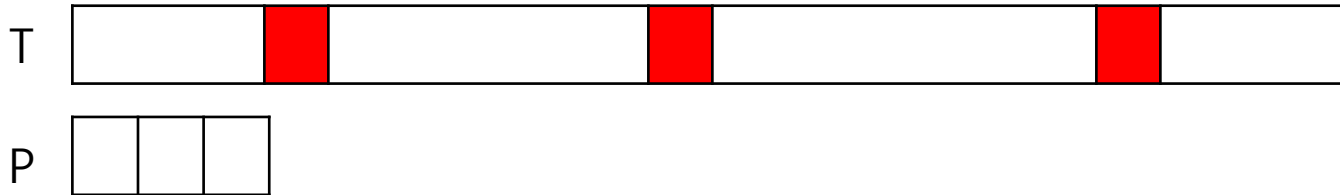
-  $r = \left\lfloor \frac{n}{k+1} \right\rfloor$



# BYP method

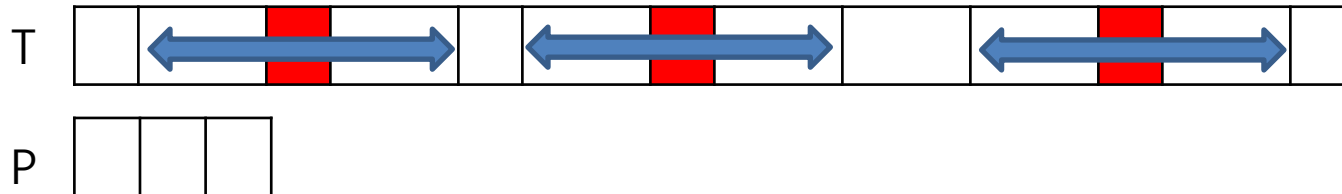
---

- Search Phase
  - Build a keyword tree for the set of "patterns"  $P$
  - Using Aho-Corasik algorithm, find  $I$ , the set of all starting locations in  $T$  where any pattern in  $P$  occurs exactly



# BYP method

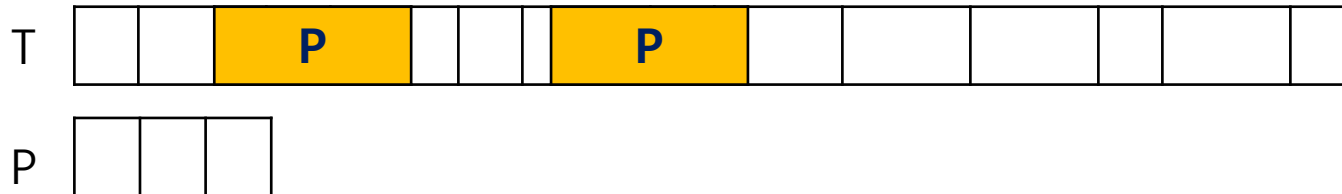
- Check Phase
  - For each index  $i \in I$  use an approximate matching algorithm to locate the end points of all approximate occurrences of  $P$  in the substring  $T[i-n-k \dots i+n+k]$



# BYP method

---

- Check Phase
  - For each index  $i \in I$  use an approximate matching algorithm to locate the end points of all approximate occurrences of  $P$  in the substring  $T[i-n-k \dots i+n+k]$





- **Time Complexity**

- a. Let  $P$  be the set of  $k+1$  substrings of  $P$  taken from the first  $k+1$  regions of  $P$ 's partition.

- b. Build a keyword tree for the set of "patterns"  $P$

- $O(n)$

- c. Using Aho-Corasik algorithm

- $O(m)$

- d. Approximate matching algorithm

- $O(m)$  – no spaces allowed, know exactly where in  $P$  any match found

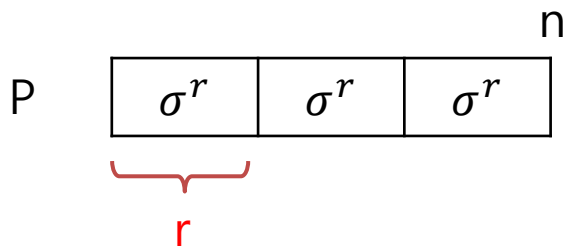
- **Expected Time**

- $\sigma$ : alphabet size
- $r$ : partition length
- any pattern  $p \in P$

# BYP method

- **Expected Time**

- $\sigma$ : alphabet size
- $r$ : partition length
- any pattern  $p \in P$ 
  - $p$  has length  $r$
  - Expected arbitrary string =  $\sigma^r$



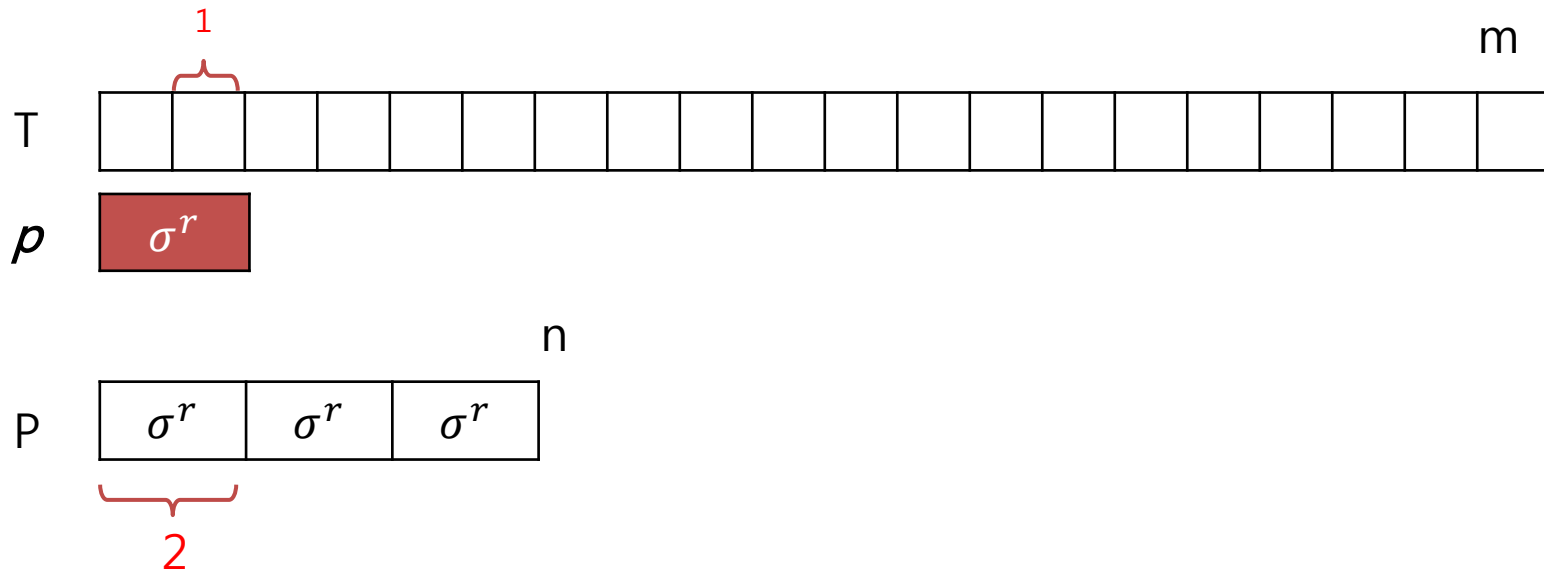
# BYP method

- Expected Time

- T contains roughly  $m$  substrings

- Expected number of exact occurrences of  $p$  in T :  $\frac{m}{\sigma^r}$

- If  $k+1 = 3$  and  $n = 6$ ,  $r = 2$



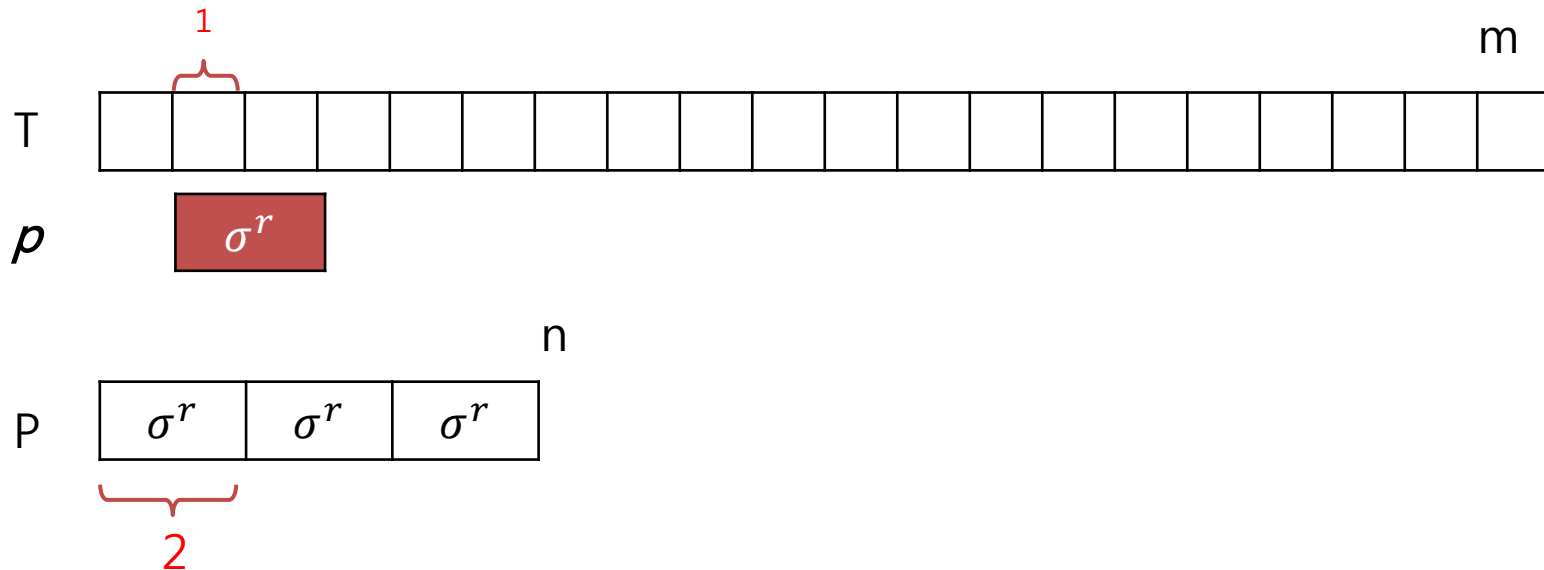
# BYP method

- Expected Time

- T contains roughly  $m$  substrings

- Expected number of exact occurrences of  $p$  in T  $= \frac{m}{\sigma^r}$

- If  $k+1 = 3$  and  $n = 6$ ,  $r = 2$



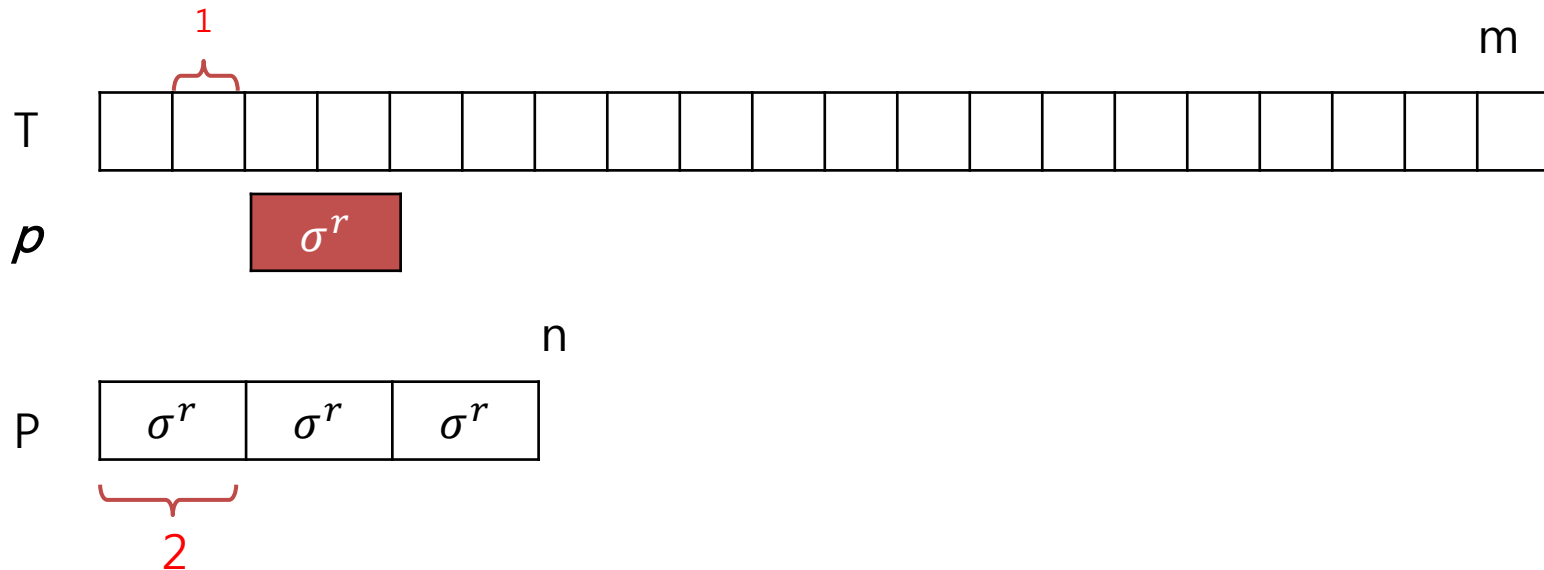
# BYP method

- Expected Time

- T contains roughly  $m$  substrings

- Expected number of exact occurrences of  $p$  in T  $= \frac{m}{\sigma^r}$

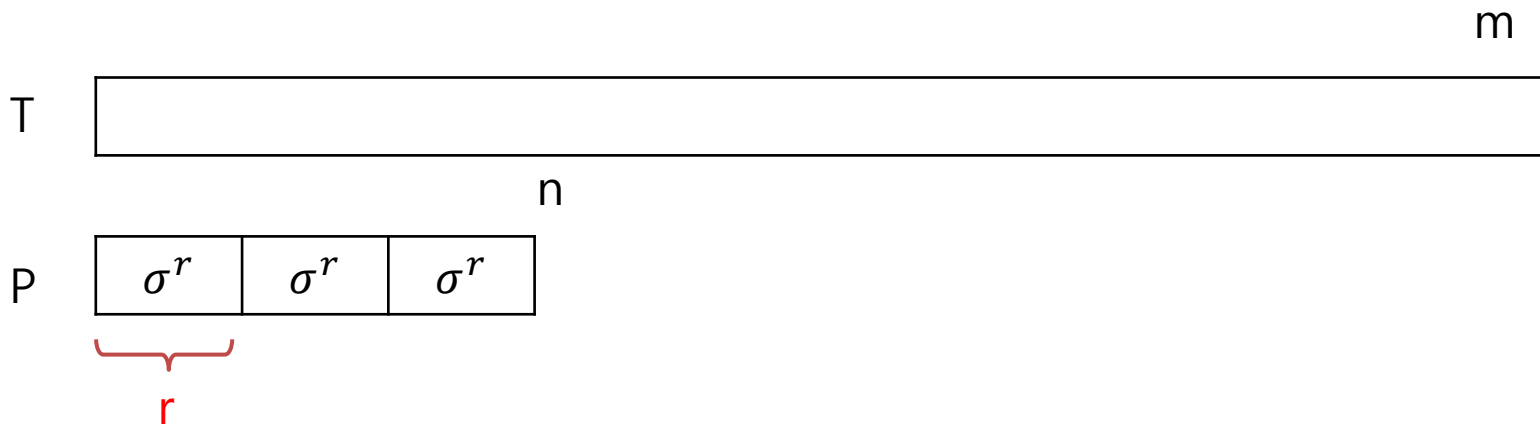
- If  $k+1 = 3$  and  $n = 6$ ,  $r = 2$



# BYP method

- **Expected Time**

- There are  $k+1$  number of pattern  $p$  in set  $P$
- Expected total number of occurrences in  $T$  of patterns from  $P$ :  $\frac{m(k+1)}{\sigma^r}$



- **Expected Time**

- Expected size of  $\mathcal{I}$ (set of survival intervals) :  $\frac{m(k+1)}{\sigma^r}$

- Expected Checking time [phase algorithm( $O(n^2)$ )] :  $\frac{mn^2(k+1)}{\sigma^r}$

- **Goal !**

- make the expected checking time linear in  $m$  for modest  $k$

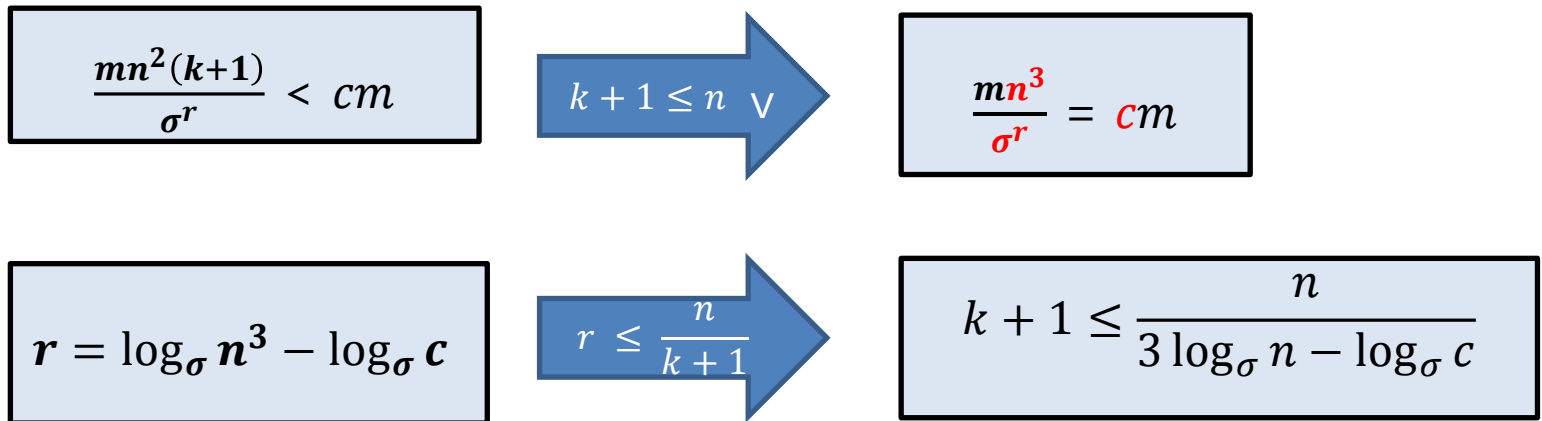
- $\frac{mn^2(k+1)}{\sigma^r} < cm$



## BYP method

- Expected Time

- make the expected checking time linear in  $m$  for modest  $k$



➤ Algorithm BYP runs in  $O(m)$  time for  $k = O(\frac{n}{\log n})$

- **Expected Time**

- $\sigma$ : alphabet size
- $r$ : partition length
- any pattern  $p \in P$

➤ Expected number of exact occurrences of  $p$  in  $T$  :  $\frac{m}{\sigma^r}$

➤ Expected total number of occurrences in  $T$  of patterns from  $P$  :  $\frac{m(k+1)}{\sigma^r}$

➤ Expected Checking time [phase algorithm( $O(n^2)$ )] :  $\frac{mn^2(k+1)}{\sigma^r}$

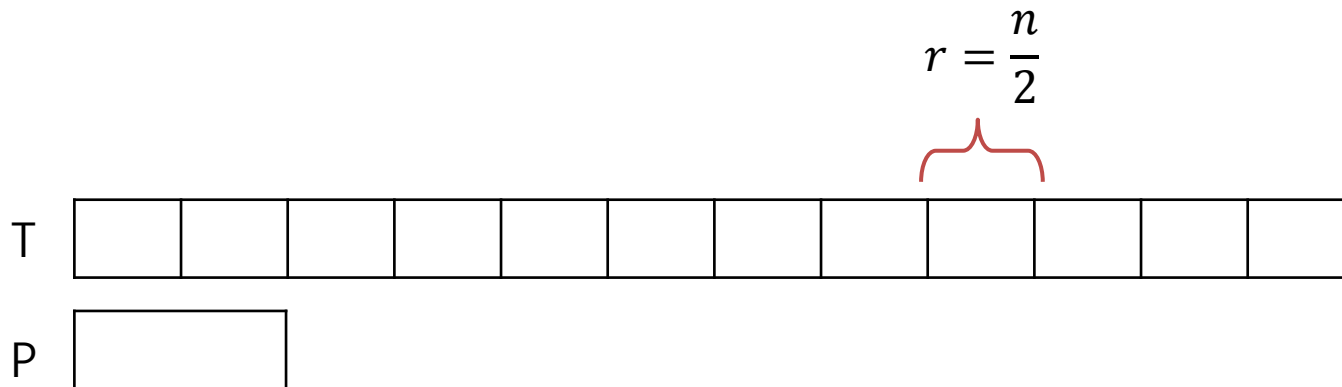
➤ Algorithm BYP runs in  $O(m)$  time for  $k = O(\frac{n}{\log n})$

- Chang-Lawler method
- $O(m)$  for  $k < n/\log_o n$
- Partition  $r = \frac{n}{2}$

## CL method

- Partition

-string T is partitioned into consecutive fixed regions of length  $r = \frac{n}{2}$

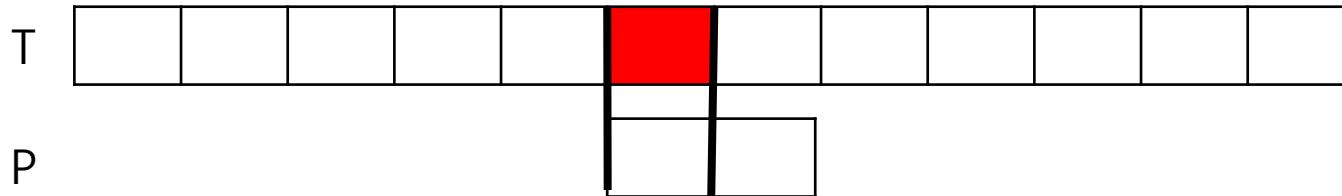


## CL method

---

- Partition

-string T is partitioned into consecutive fixed regions of length  $r = \frac{n}{2}$

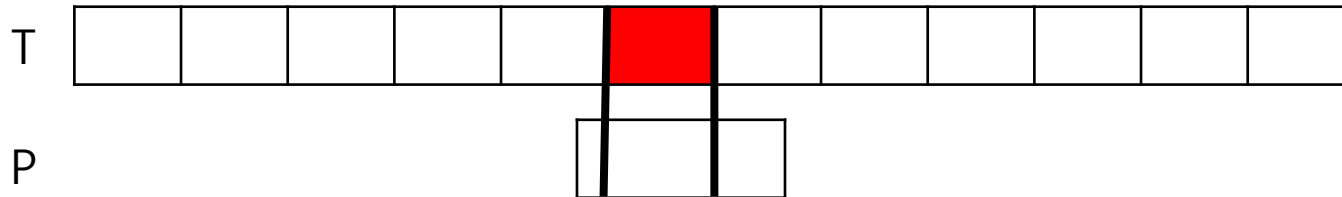


## CL method

---

- Partition

-string T is partitioned into consecutive fixed regions of length  $r = \frac{n}{2}$

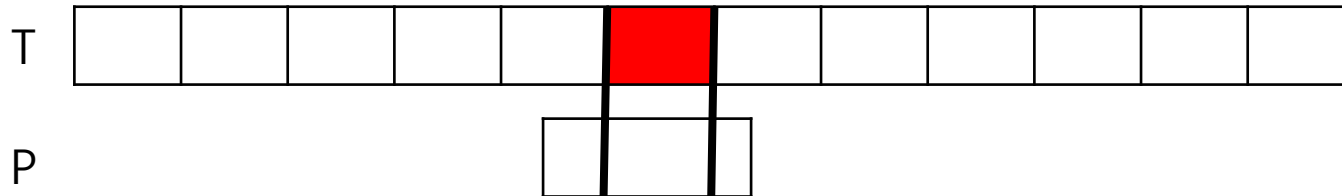


## CL method

---

- Partition

-string T is partitioned into consecutive fixed regions of length  $r = \frac{n}{2}$

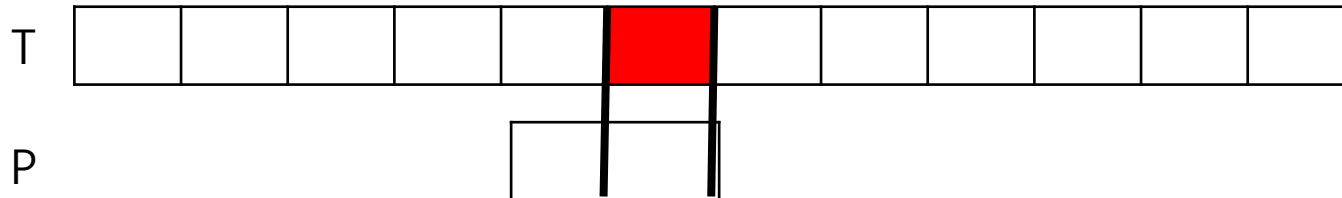


## CL method

---

- Partition



-string T is partitioned into consecutive fixed regions of length  $r = \frac{n}{2}$

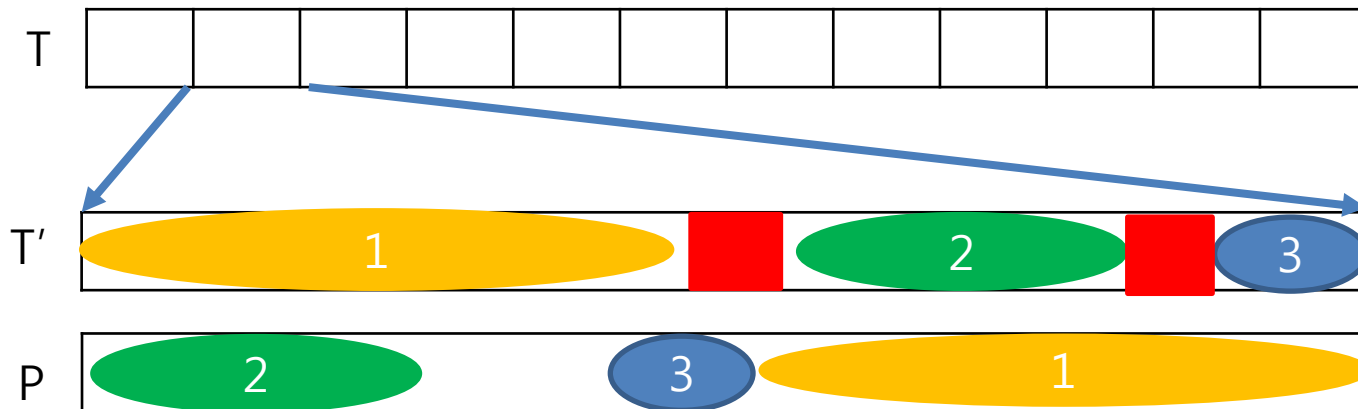




## CL method

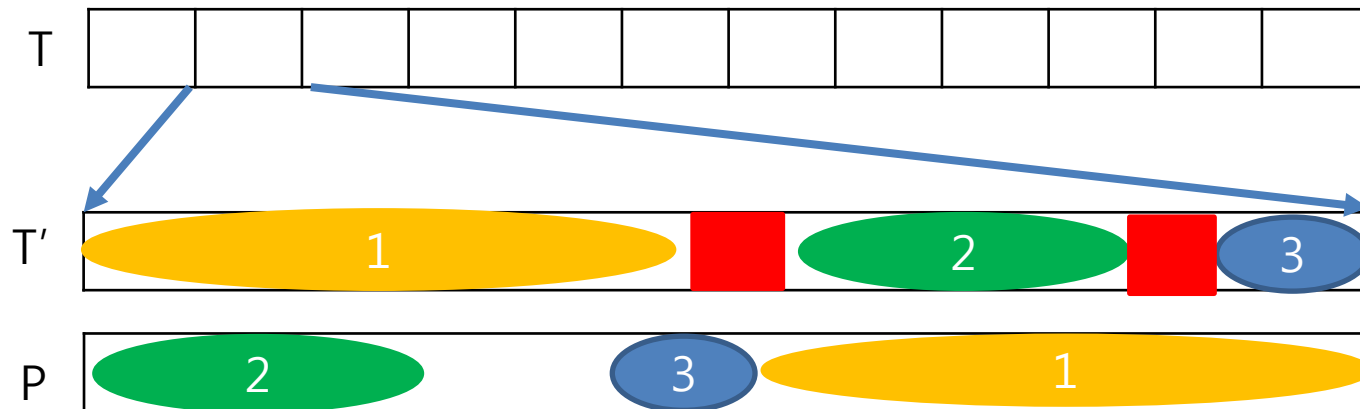
- Search Phase
  - Use *matching statistics*
  - $ms(i)$  : length of the longest substring starting at position  $i$  of  $T$  that matches a substring *somewhere* in  $P$

 Longest substring  
 mismatch



## CL method

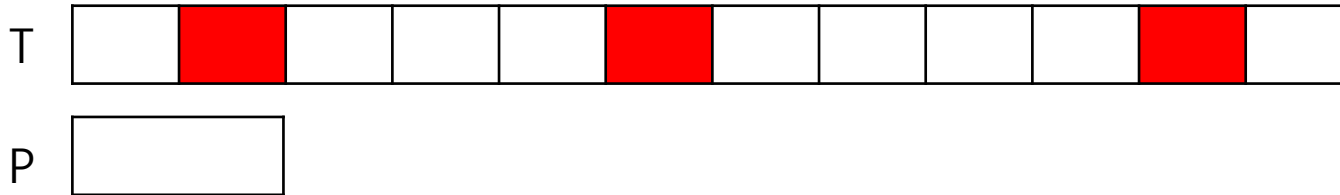
- Search Phase
  - If mismatch exceed  $k$ ,  
stop searching and exclude region from set  $R$  (surviving region)
  - If  $T'$  has finished being mapped,  
leave the region in set  $R$  and go onto the next region.



## CL method

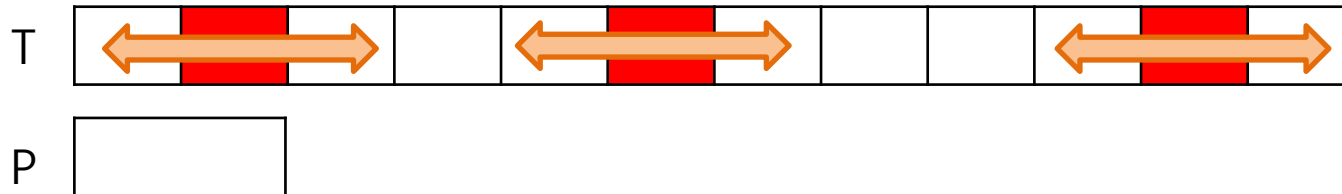
---

- Search Phase
  - If mismatch exceed  $k$ ,  
stop searching and exclude region from set  $R$ (surviving region)
  - If  $T'$  has finished being mapped,  
leave the region in set  $R$  and go onto the next region.



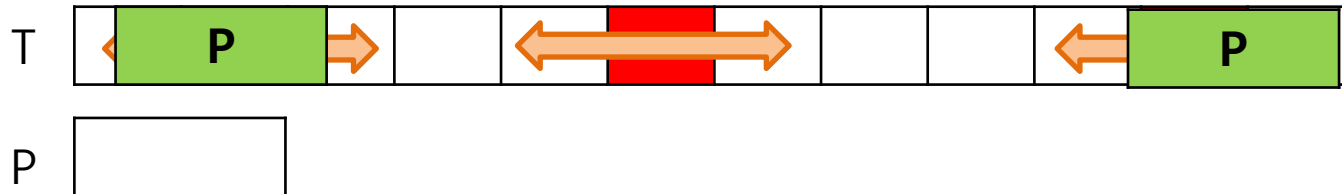
## CL method

- Check Phase
  - R is a surviving region
  - Executes approximate matching algorithm for P against a neighborhood of T that starts  $n/2$  positions to the left of R and ends  $n/2$  positions to its right



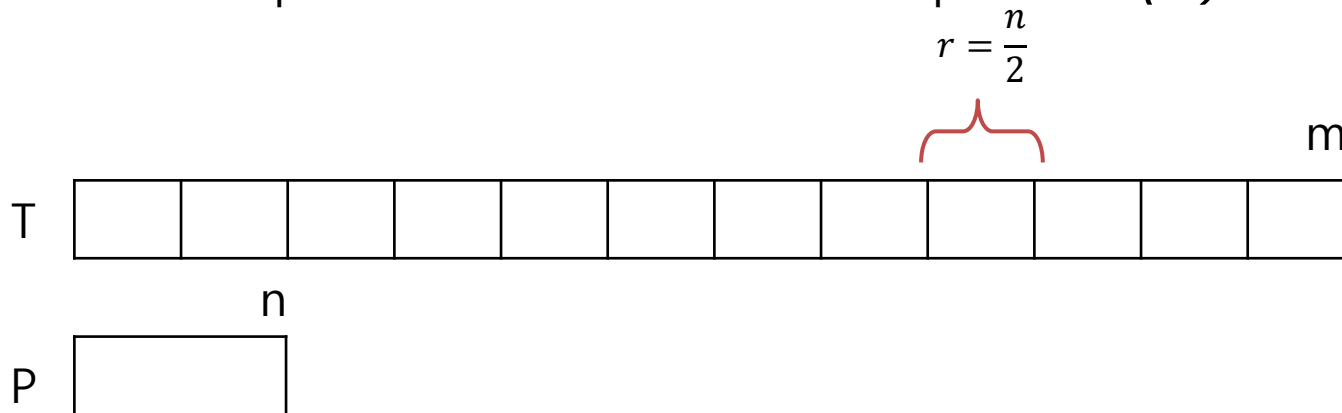
## CL method

- Check Phase
  - R is a surviving region
  - Executes approximate matching algorithm for P against a neighborhood of T that starts  $n/2$  positions to the left of R and ends  $n/2$  positions to its right



# CL method

- Time Complexity
    - Search phase
      - number of regions =  $m(n/2) = 2m/n$
      - $ms(i)$  is repeated until either mismatches exceed  $k$  or the whole region was searched.
      - $E(M)$  : expected value of matching statistics.
- => Expected value is less than or equal to  $kE(M)$



## CL method – search phase

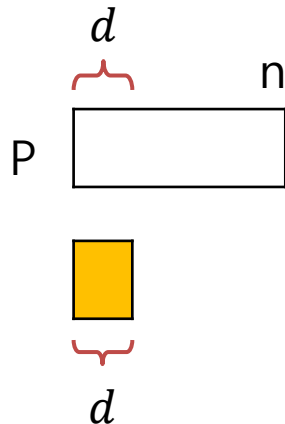
---

- **Lemma 12.3.3**

$E(M)$ , the expected value of a matching statistic, is  $O(\log_\sigma n)$

- **Proof>**

- For fixed length  $d$ , roughly  $n$  substrings of length  $d$  in  $P$



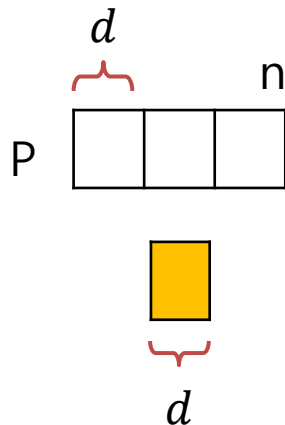
## CL method – search phase

- **Lemma 12.3.3**

$E(M)$ , the expected value of a matching statistic, is  $O(\log_\sigma n)$

- **Proof>**

- For fixed length  $d$ , roughly  $n$  substrings of length  $d$  in  $P$





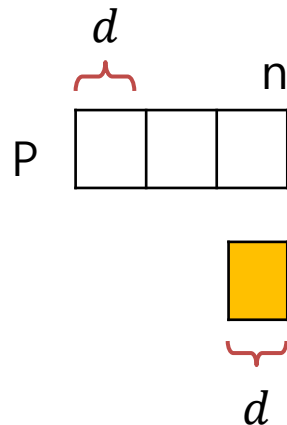
## CL method – search phase

- **Lemma 12.3.3**

$E(M)$ , the expected value of a matching statistic, is  $O(\log_\sigma n)$

- **Proof>**

- For fixed length  $d$ , roughly  $n$  substrings of length  $d$  in  $P$



## CL method – search phase

---

- **Lemma 12.3.3**

$E(M)$ , the expected value of a matching statistic, is  $O(\log_{\sigma} n)$

- **Proof>**

- For fixed length  $d$ , roughly  $n$  substrings of length  $d$  in  $P$
- For any specific string  $\alpha$  of length  $d$ ,  
probability that  $\alpha$  is found somewhere in  $P < \frac{n}{\sigma^d}$
- unless  $\sigma^d = n$

## CL method – search phase

---

- **Lemma 12.3.3**

$E(M)$ , the expected value of a matching statistic, is  $O(\log_\sigma n)$

- **Proof>**

- Let  $X$  be the random variable that has value  $\log_\sigma n$  for  $ms(i) \leq \log_\sigma n$
- Fixed  $d$  until infinite
- Then  $E(M) = O(\log_\sigma n)$

$$E(M) < E(X) < \log_\sigma n + \sum_{l=\log_\sigma n}^{\infty} \frac{l}{\sigma^l} = \log_\sigma n + 2.$$

## CL method

---

- Time Complexity
  - Search phase
    - number of regions =  $m(n/2) = 2m/n$
    - Matching statistics:  $kE(M) = O(k \log_{\sigma} n)$
  - Total expected search phase
    - =  $O(2mkE(M)/n)$
    - =  $O(2mk \log_{\sigma} n/n)$
    - same process as the BYP method
    - =  **$O(m)$  for  $k < n/\log_{\sigma} n$**

- Time Complexity
  - Check Phase
    - Expected number of regions that survived :  $e$
    - Total check phase time  
 $= O(kne)$   
*too difficult to present analysis of  $e$*
    - $= O(km/n^3)$  when  $e = m/n^4$   
 $= \mathbf{O(m)}$  when  $k = O(n/\log_{\sigma} n)$ ,

## Problems with previous methods

---

1. They permit a large expected number of surviving regions compared to the expected number of true approximate matches
2. When a surviving region is first located, the method move directly to full dynamic programming computations

## Myers method

---

- Handle insertion, deletion as well as mismatches
- $O(km^{p(\epsilon)} \log m)$  when length of the sub pattern is  $\log_o n$
- Too complex for detailed discussion

# Myers method

## • Lemma 12.3.4

- Suppose  $\alpha$   $\varepsilon$ -matches  $\beta$ . Then  $\beta$  can be divided into two substrings  $\beta_0$  and  $\beta_1$  such that  $\beta = \beta_0\beta_1$ , and either  $\alpha_0$   $\varepsilon$ -matches  $\beta_0$  or  $\alpha_1$   $\varepsilon$ -matches  $\beta_1$ .
- Let  $\alpha = \alpha_0\alpha_1$ , where  $|\alpha_0|$  is assumed equal to  $|\alpha_1|$ .

$$\text{Error rate} : \frac{8}{24} = \frac{1}{3}$$

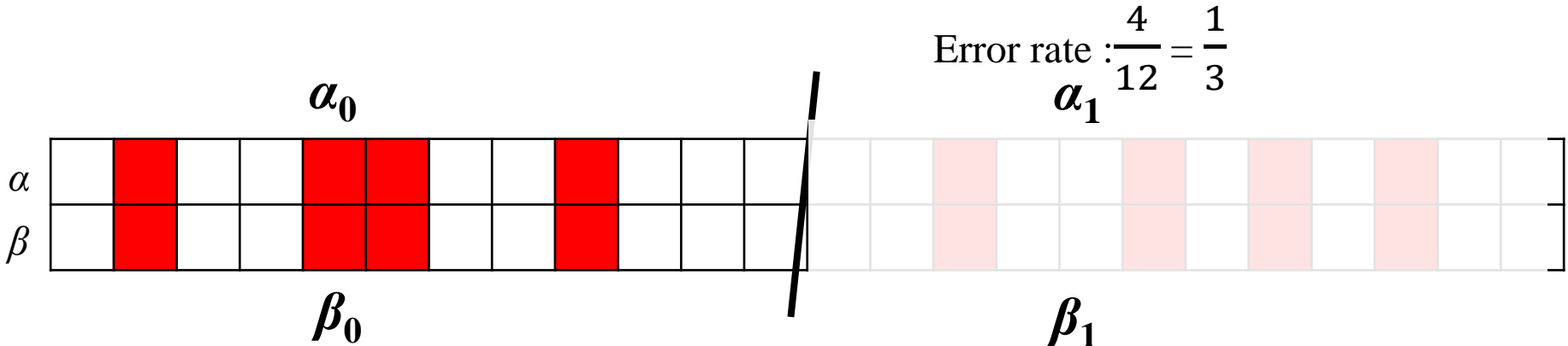
$\alpha$																							
$\beta$																							



© 2004 Blackwell Publishing Ltd, *Journal of Internal Medicine* 255: 105–112

- **Lemma 12.3.4**

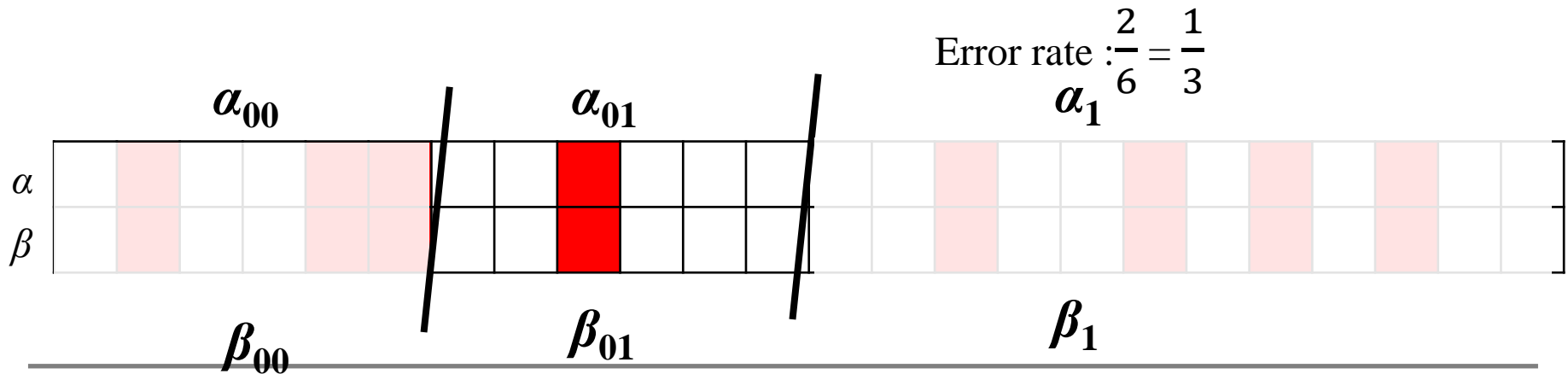
- Suppose  $\alpha$   $\varepsilon$ -matches  $\beta$ . Then  $\beta$  can be divided into two substrings  $\beta_0$  and  $\beta_1$  such that  $\beta = \beta_0\beta_1$ , and either  $\alpha_0$   $\varepsilon$ -matches  $\beta_0$  or  $\alpha_1$   $\varepsilon$ -matches  $\beta_1$ .
- Let  $\alpha = \alpha_0\alpha_1$ , where  $|\alpha_0|$  is assumed equal to  $|\alpha_1|$ .



# Myers method

## • Lemma 12.3.4

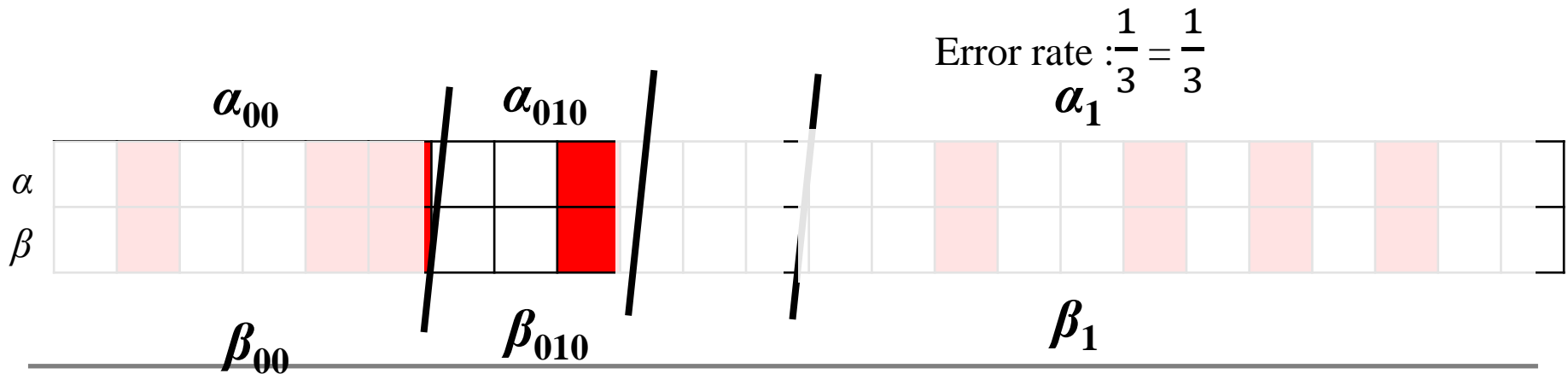
- Suppose  $\alpha$   $\varepsilon$ -matches  $\beta$ . Then  $\beta$  can be divided into two substrings  $\beta_0$  and  $\beta_1$  such that  $\beta = \beta_0\beta_1$ , and either  $\alpha_0$   $\varepsilon$ -matches  $\beta_0$  or  $\alpha_1$   $\varepsilon$ -matches  $\beta_1$ .
- Let  $\alpha = \alpha_0\alpha_1$ , where  $|\alpha_0|$  is assumed equal to  $|\alpha_1|$ .



# Myers method

## • Lemma 12.3.4

- Suppose  $\alpha$   $\varepsilon$ -matches  $\beta$ . Then  $\beta$  can be divided into two substrings  $\beta_0$  and  $\beta_1$  such that  $\beta = \beta_0\beta_1$ , and either  $\alpha_0$   $\varepsilon$ -matches  $\beta_0$  or  $\alpha_1$   $\varepsilon$ -matches  $\beta_1$ .
- Let  $\alpha = \alpha_0\alpha_1$ , where  $|\alpha_0|$  is assumed equal to  $|\alpha_1|$ .



# Myers method

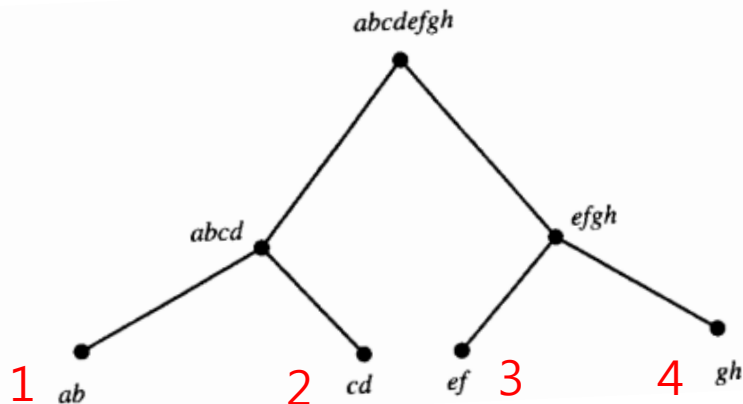
---

- d-neighborhood
  - example>
  - only {a,b}, S = aba, d=1
  - 1-neighborhood : {bba,aaa,abb,aaba,abaa,baba,abba,abab,ba,aa,ab}
  - Condense(remove prefix): {bba, aaa, aaba, abaa, baba, abba, abab}

# Myers method

- First iteration
  - T has already been preprocessed into some index structure
  - Partition P into subpatterns of length  $\log_{\sigma} m$
  - Construct the condensed d-neighborhood for each subpattern in  $P$
  - Find all locations of substrings in test T that exactly match one of the substrings in one of the condensed d- neighborhoods.

Binary tree of P



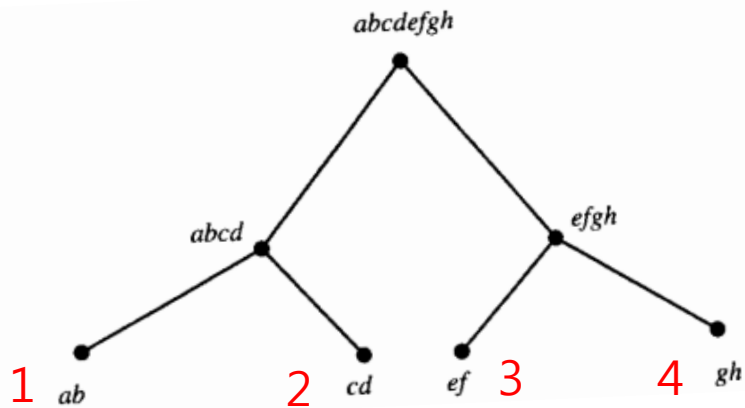
T



# Myers method

- Successive iteration
  - extend each initial surviving match to become a  $\varepsilon$ -match between substring twice as long as those in the current surviving match

Binary tree of P



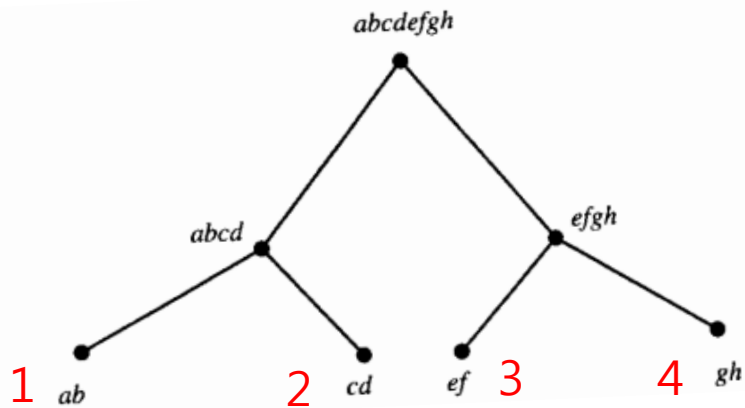
T

1	2					3	4					3	4			1	2				
---	---	--	--	--	--	---	---	--	--	--	--	---	---	--	--	---	---	--	--	--	--

# Myers method

- Successive iteration
  - extend each initial surviving match to become a  $\varepsilon$ -match between substring twice as long as those in the current surviving match

Binary tree of P



T

1	2	3	4	1	2	3	4			1	2	3	4			1	2	3	4		
---	---	---	---	---	---	---	---	--	--	---	---	---	---	--	--	---	---	---	---	--	--

# Myers method

---

- Time Complexity
  - when length of the sub pattern is  $\log_{\sigma} n$
  - total runtime is  $O(km^{p(\epsilon)} \log m)$
  - $p(\epsilon)$  convex increases very slowly



## Myers method

---

- Exposition given above is only the outline of Myer's method, without any analysis
- Unlike BYP and CL, error rates that establish sublinear running times do not depend on the length of  $P$
- Although expected running times for both CL and Myers' are sublinear
  - CL : due to multiplicative factor that is less than one
  - Myers : due to an exponent that is less than one

## Final comment of exclusion method

---

- Developed with motivation of searching large DNA and protein databases for approximate occurrences of query strings
- weak for the case of **protein database search**, error rates as high as 85% are great interest when comparing protein sequences
- The most effective practical database search methods can be considered as exclusion methods  
ex ) BLAST, FASTA, and variants