

BWT Arrays and Mismatching Trees: A New Way for String Matching with k Mismatches

Yangjun Chen, Yujia Wu
Department of Applied Computer Science
University of Winnipeg

Outline

- **Motivation**
 - Statement of Problem
 - Related work
- **BWT Arrays – A space-economic Index for String Matching**
- **String Matching with k Mismatches**
 - Search trees
 - Mismatching information
 - Mismatching trees
- **Experiments**
- **Conclusion and Future Work**

Statement of Problem

- **String matching with k mismatches:** find all the occurrences of a pattern string r in a target string s with each occurrence having up to k positions different between r and s .
 - In DNA databases, due to polymorphisms or mutations among individuals or even sequencing errors, a read (a short sample DNA sequence) may disagree in some positions at any of its occurrences in a genome.

Example: $k = 4$

a	a	a	a	a	c	a	a	a	c	←-----	pattern				
a	c	a	c	a	c	a	g	a	a	g	c	c	c	←-----	target

Related Work

➤ Exact string matching

- On-line algorithms:

Knuth-Morris-Pratt, Boyer-Moore, Aho-Corasick

- Index based:

suffix trees (Weiner; McCreight; Ukkonen), suffix arrays (Manber, Myers), BWT-transformation (Burrow-Wheeler), Hash (Karp, Rabin)

➤ Inexact string matching

- String matching with k mismatches - Hamming distance (*Landau, U. Vishkin; Amir et al.; Cole*)

- String matching with k differences - Levenshtein distance (*Chang, Lampe*)

- String matching with wild-cards (*Manber, Baeza-Yates*)

BWT-Index

➤ Burrows-Wheeler Transform (BWT)

➤ $s = a_1c_1a_2g_1a_3c_2a_4\$$

Rank correspondence:

	\underline{rk}_F	F	L	\underline{rk}_L
$a_1 c_1 a_2 g_1 a_3 c_2 a_4 \$$	-	$\$$	$a_1 c_1 a_2 g_1 a_3 c_2 a_4$	1
$c_1 a_2 g_1 a_3 c_2 a_4 \$ a_1$	1	a_4	$\$ a_1 c_1 a_2 g_1 a_3 c_2$	1
$a_2 g_1 a_3 c_2 a_4 \$ a_1 c_1$	2	a_3	$c_2 a_4 \$ a_1 c_1 a_2 g_1$	1
$g_1 a_3 c_2 a_4 \$ a_1 c_1 a_2$	3	a_1	$c_1 a_2 g_1 a_3 c_2 a_4 \$$	-
$a_3 c_2 a_4 \$ a_1 c_1 a_2 g_1$	4	a_2	$g_1 a_3 c_2 a_4 \$ a_1 c_1$	2
$c_2 a_4 \$ a_1 c_1 a_2 g_1 a_3$	1	c_2	$a_4 \$ a_1 c_1 a_2 g_1 a_3$	2
$a_4 \$ a_1 c_1 a_2 g_1 a_3 c_2$	2	c_1	$a_2 g_1 a_3 c_2 a_4 \$$	3
$\$ a_1 c_1 a_2 g_1 a_3 c_2 a_4$	1	g_1	$a_3 c_2 a_4 \$ a_1 c_1 a_2$	4

rank: 3

rank: 3

$$\underline{rk}_F(e) = \underline{rk}_L(e)$$

BWT construction:

$$\begin{cases} L[i] = \$, & \text{if } SA[i] = 1; \\ L[i] = s[SA[i] - 1], & \text{otherwise.} \end{cases}$$

$SA[\dots]$ – suffix array

Backward Search of BWT-Index

- $s = a_1c_1a_2g_1a_3c_2a_4\$$
- Search $p = aca$

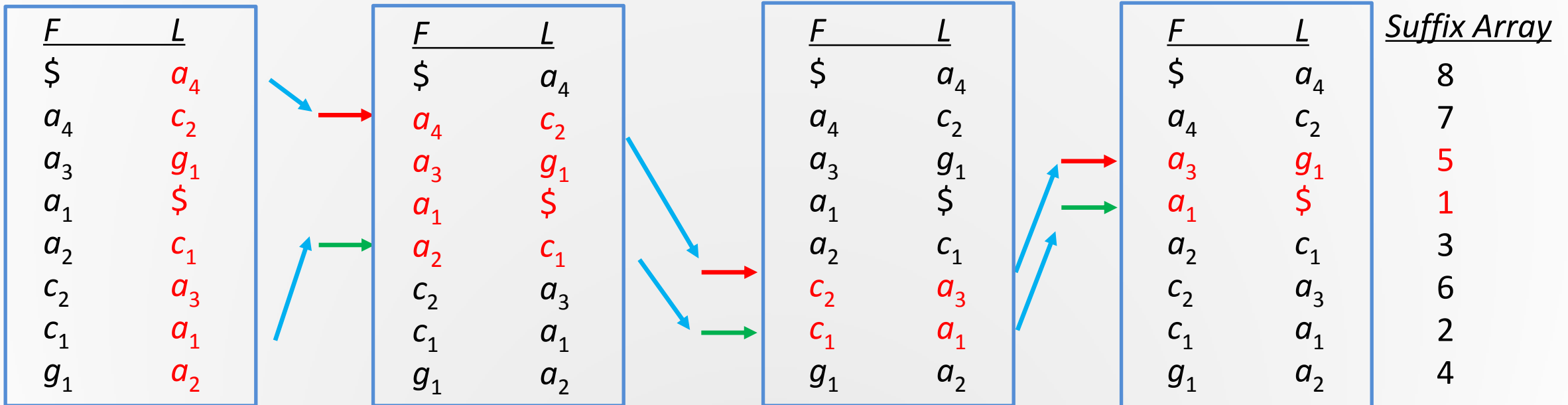
←--- Backward Search

$$\text{search}(z, \pi) = \begin{cases} \langle z, [\alpha, \beta] \rangle, & \text{if } z \text{ appears in } L_\pi; \\ \phi, & \text{otherwise.} \end{cases}$$

z : a character

π : a range in F

L_π : a range in L , corresponding to π



Backward Search of BWT-Index

$\text{search}(c, \langle a, [2, 5] \rangle)$

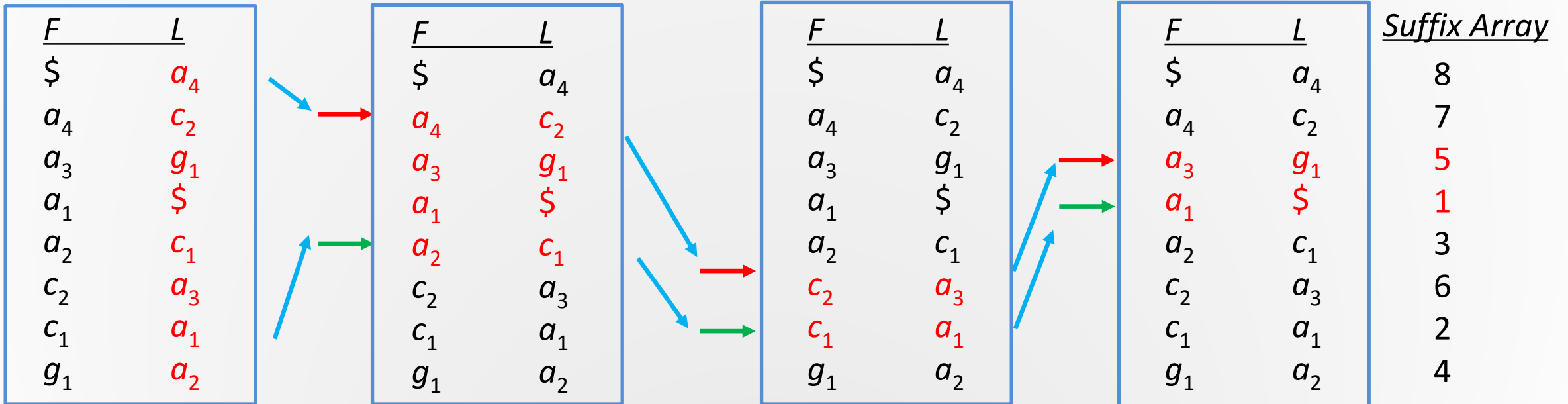
$\text{search}(a, \langle c, [1, 2] \rangle)$

Search sequence:

$\langle a, [2, 5] \rangle$

$\langle c, [1, 2] \rangle$

$\langle a, [3, 4] \rangle$



rankAll

- Arrange $|\Sigma|$ arrays each for a character $x \in \Sigma$ such that $A_x[i]$ (the i th entry in the array for x) is the number of appearances of x within $L[1 \dots i]$.
- Instead of scanning a certain segment $L[\alpha \dots \beta]$ ($\alpha \leq \beta$) to find a subrange for a certain $x \in \Sigma$, we can simply look up A_x to see whether $A_x[\alpha - 1] = A_x[\beta]$. If it is the case, then α does not occur in $\alpha \dots \beta$. Otherwise, $[A_x[\alpha - 1] + 1, A_x[\beta]]$ should be the found range.

Example

To find the first and the last appearance of c in $L[2 \dots 5]$, we only need to find $A_c[2 - 1] = A_c[1] = 0$ and $A_c[5] = 2$. So the corresponding range is $[A_c[2 - 1] + 1, A_c[5]] = [1, 2]$.

F	L
\$	a_4
a_4	c_2
a_3	g_1
a_1	\$
a_2	c_1
c_2	a_3
c_1	a_1
g_1	a_2

$A_\$$	A_a	A_c	A_g	A_t
0	1	0	0	0
0	1	1	0	0
0	1	1	1	0
1	1	1	1	0
1	1	2	1	0
1	2	2	1	0
1	3	2	1	0
1	4	2	1	0

Reduce *rankAll*-Index Size

- **F-ranks:** $F_\alpha = \langle \alpha; x_\alpha, y_\alpha \rangle$
- **BWT array:** L
- **Reduced appearance array:** A_α with bucket size β .
- **Reduced suffix array:** SA^* with bucket size γ .

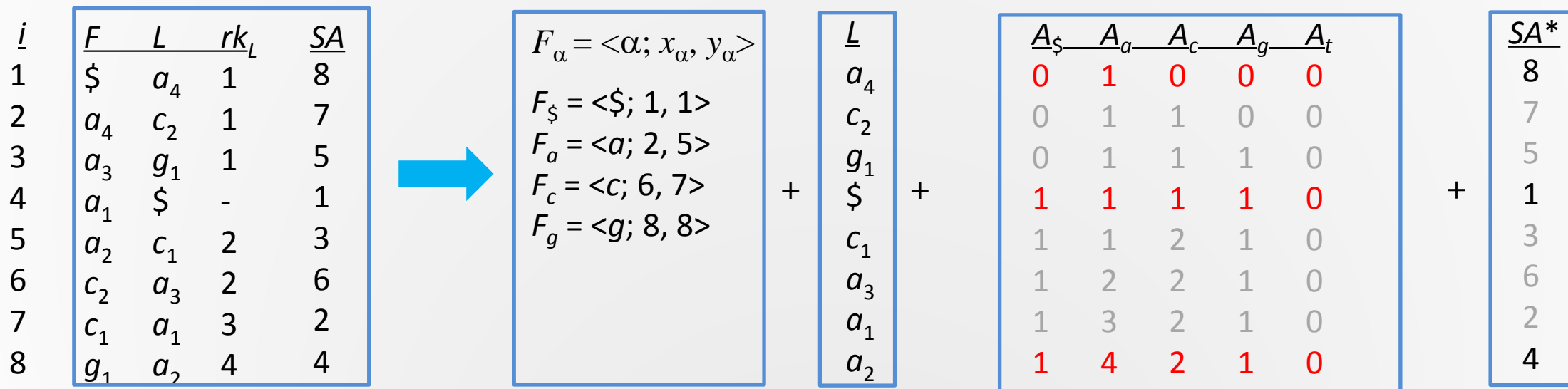
Find a range:

$$top' \leftarrow F(x_\alpha) + A_\alpha[\lfloor (top - 1)/\beta \rfloor] + r + 1$$

$$bot' \leftarrow F(x_\alpha) + A_\alpha[\lfloor bot/\beta \rfloor] + r'$$

r is the number of α 's appearances within $L[\lfloor (top - 1)/\beta \rfloor \beta .. top - 1]$

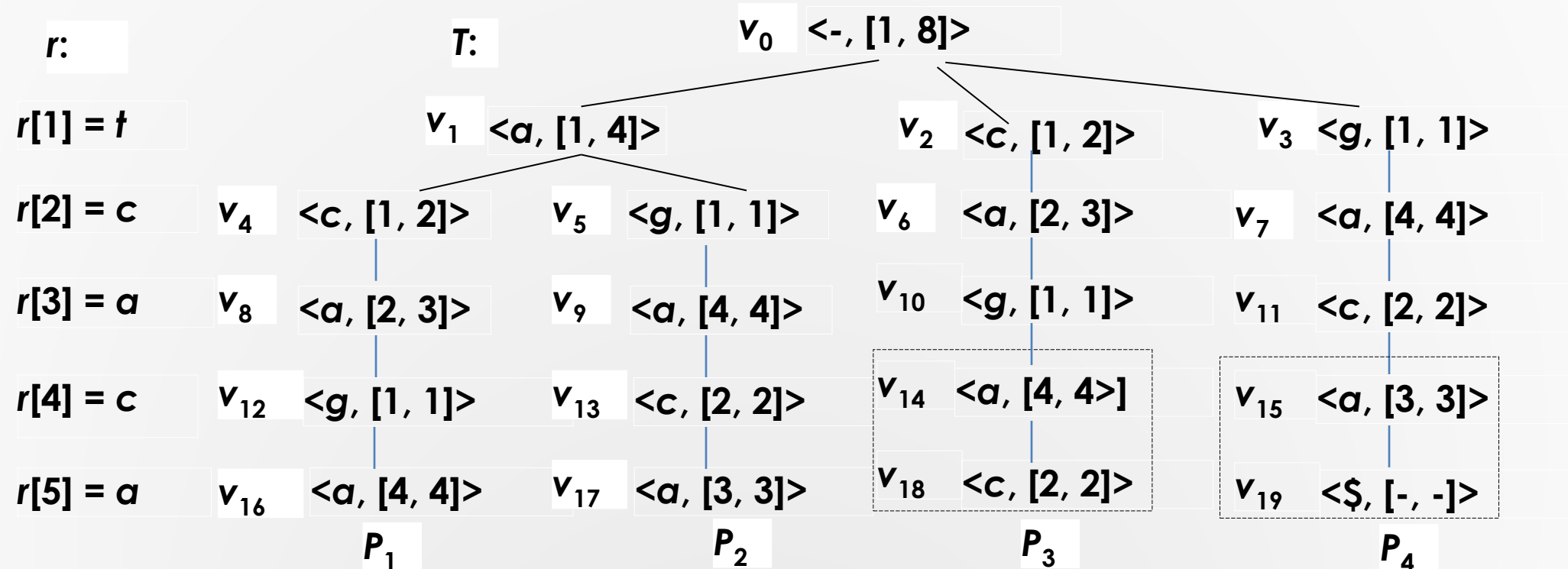
r' is the number of α 's appearances within $L[\lfloor bot/\beta \rfloor \beta .. bot]$



String Matching with k Mismatches

➤ Search Trees

pattern: $r = tcaca$; target: $s = acagaca$; $k = 2$.

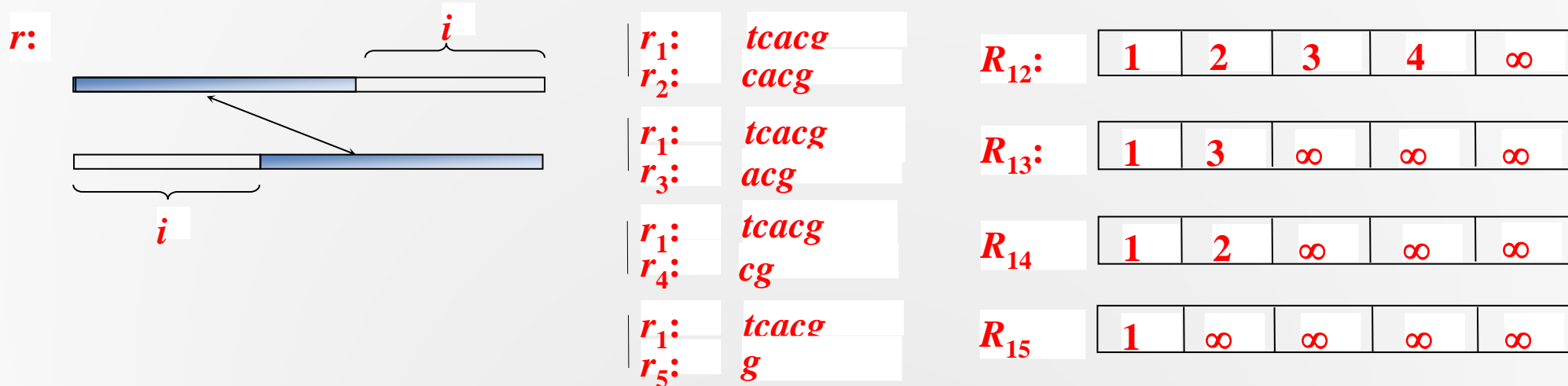


String Matching with k Mismatches

➤ Mismatching information

R – mismatching table for r with $|r| = m$.

R_{ij} – containing the positions of the first $2k + 1$ mismatches between $r[i .. m - q + i]$ and $r[j .. m - q + j]$, where $q = \max\{i, j\}$, such that if $R_{ij}[l] = x$ ($\neq \infty$) then $r[i + x - 1] \neq r[j + x - 1]$ or one of them does not exist, and it is the l -th mismatch between them.



String Matching with k Mismatches

➤ Derivation of mismatching information

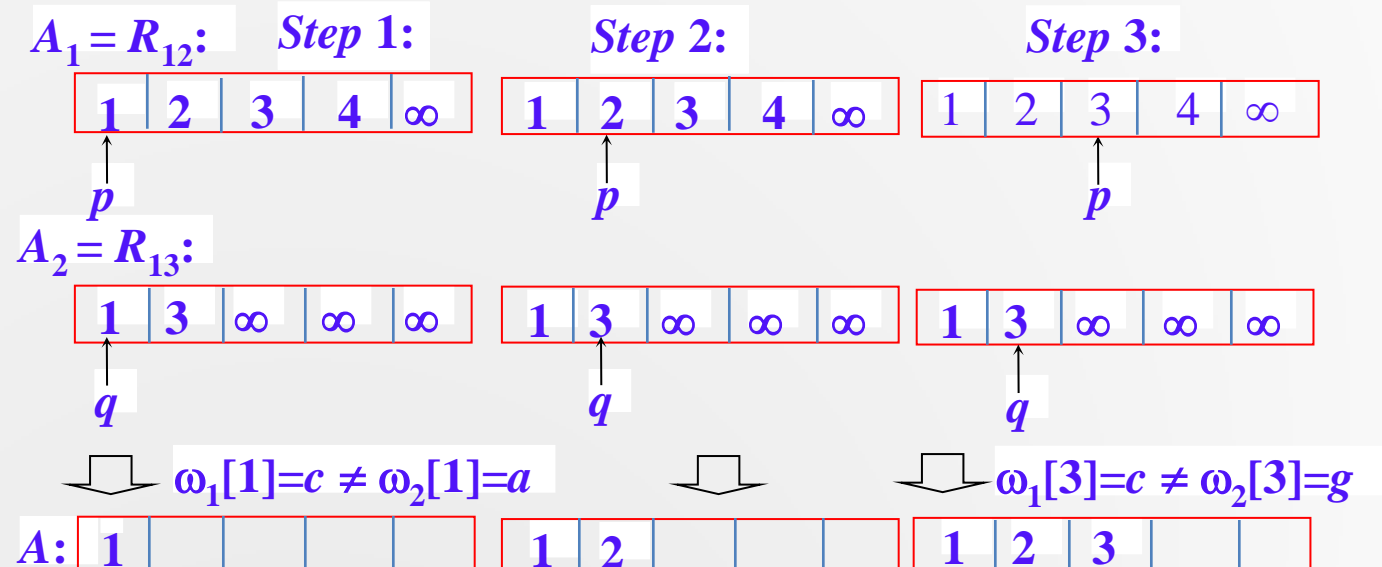
We store only part of mismatching information, specifically: R_{12}, \dots, R_{1m} , while all the other mismatching information will be dynamically derived.

Derive the mismatching information between

$\omega_1 = r[2 \dots 4] = \text{cacg}$ and

$\omega_2 = r[3 \dots 5] = \text{acg}$

from R_{12} and R_{13} .



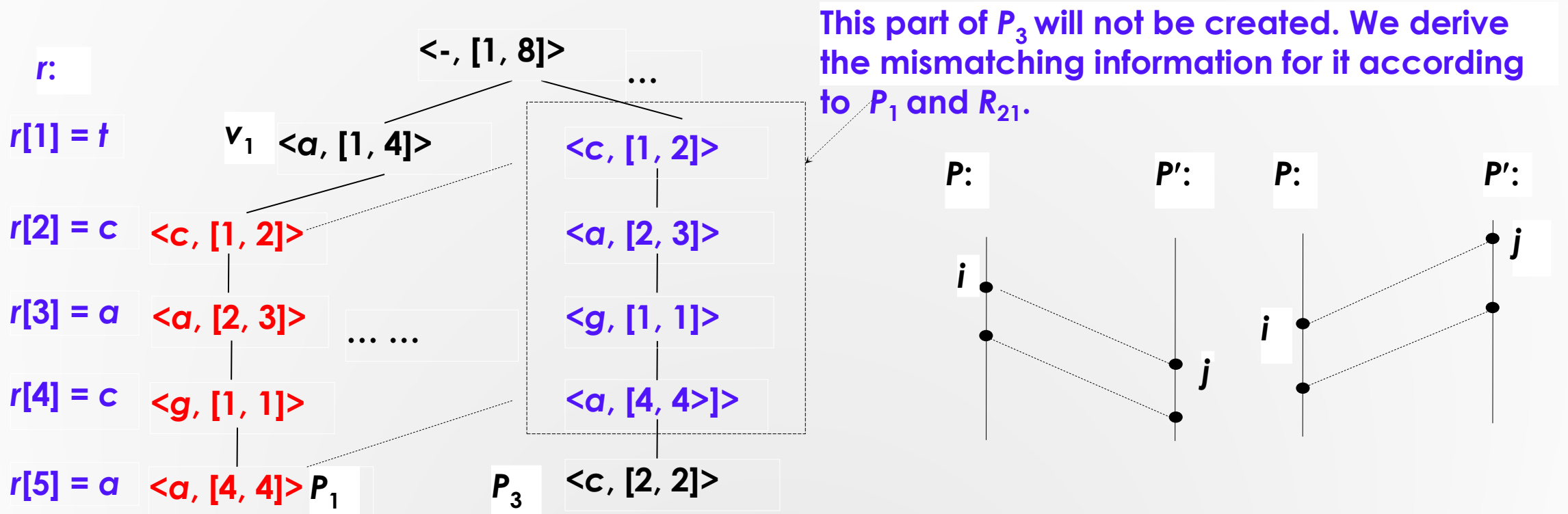
String Matching with k Mismatches

- **Algorithm for Derivation of mismatching information**
 - ❑ Let ω , ω_1 and ω_2 be three strings. Let A_1 and A_2 be two arrays containing all the positions of mismatches between ω and ω_1 , and ω and ω_2 , respectively.
 - ❑ Create a new array A such that if $A[i] = j$ ($\neq \infty$), then $\omega_1[j] \neq \omega_2[j]$, or one of them does not exist. It is the i th mismatch between them.

1. $p := 1; q := 1; l := 1;$
2. If $A_2[q] < A_1[p]$, then $\{A[l] := A_2[q]; q := q + 1; l := l + 1;\}$
3. If $A_1[p] < A_2[q]$, then $\{A[l] := A_1[p]; p := p + 1; l := l + 1;\}$
4. If $A_1[p] = A_2[q]$, then $\{\text{if } \omega_1[p] \neq \omega_2[q], \text{ then } \{A[l] := q; l := l + 1;\} p := p + 1; q := q + 1;\}$
5. If $p > |A_1|$, $q > |A_2|$, or both $A_1[p]$ and $A_2[q]$ are ∞ , stop (if A_1 (or A_2) has some remaining elements, which are not ∞ , first append them to the rear of A , and then stop.)
6. Otherwise, go to (2).

String Matching with k Mismatches

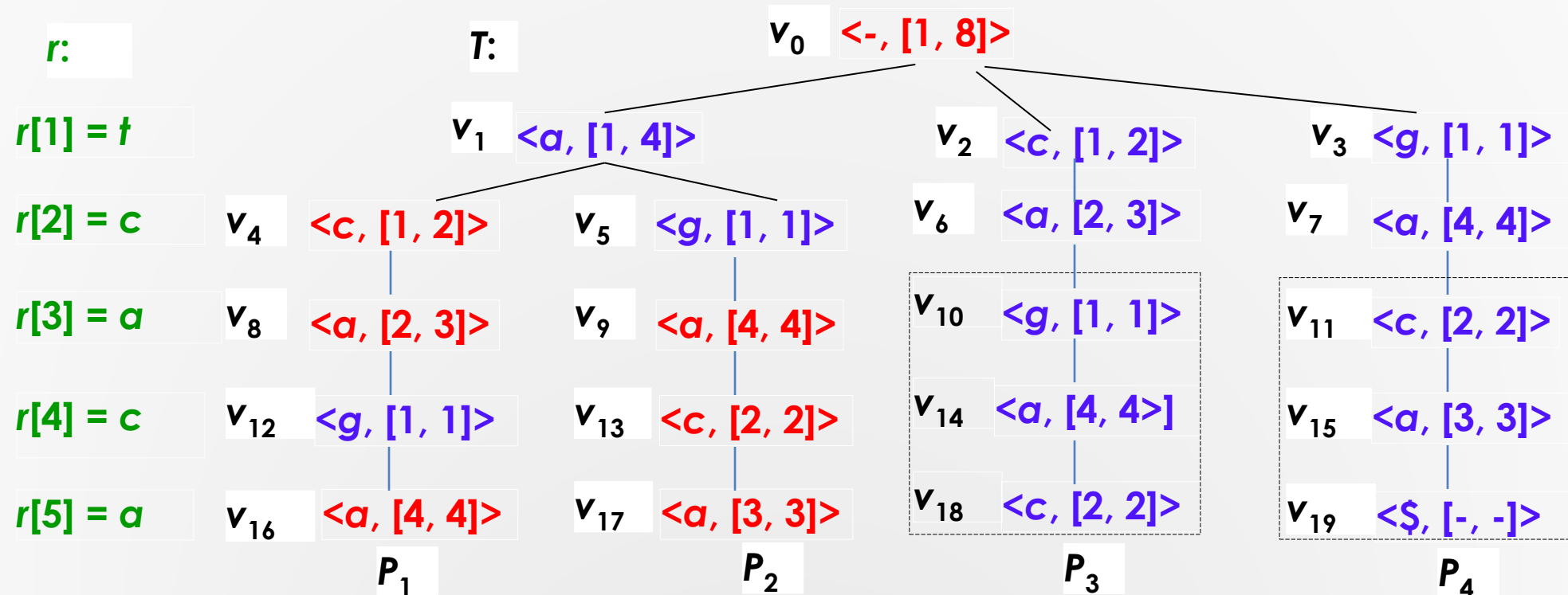
- Derivation of mismatching information for paths in a search tree.



String Matching with k Mismatches

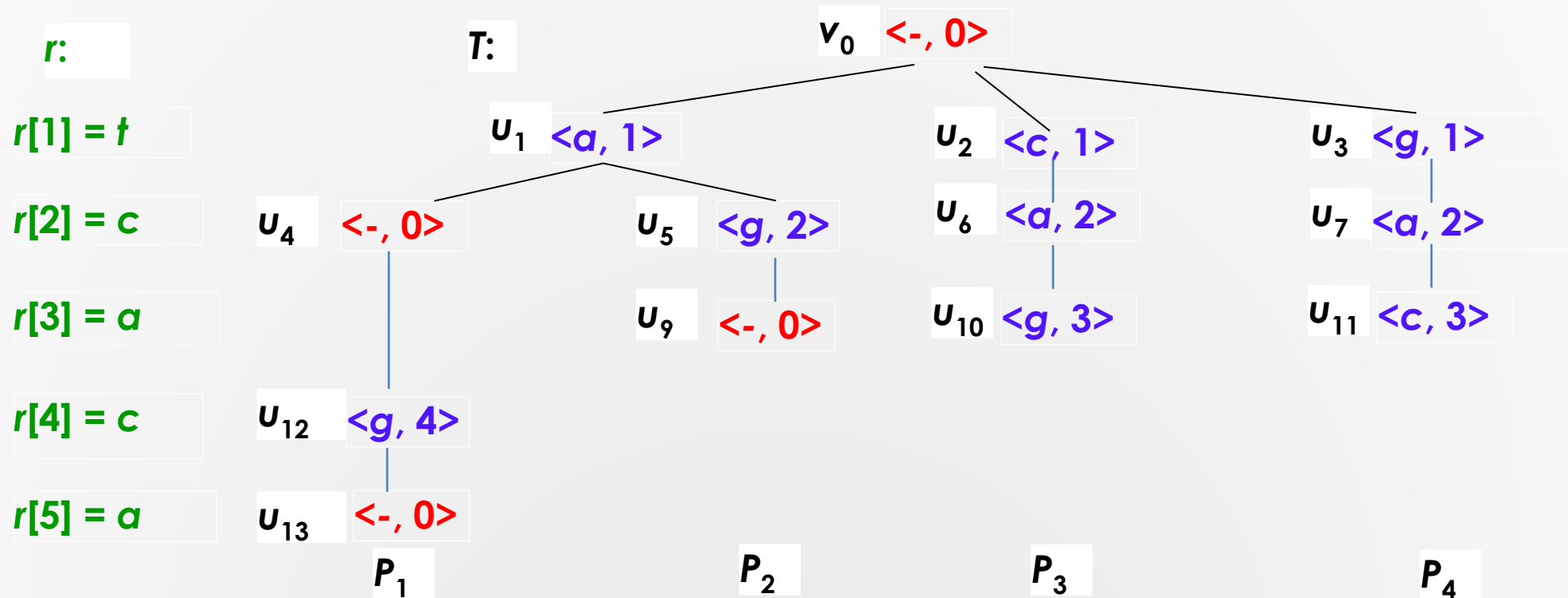
➤ Mismatching trees

In a search tree, we distinguish between matching and mismatching nodes.



String Matching with k Mismatches

➤ Mismatching trees

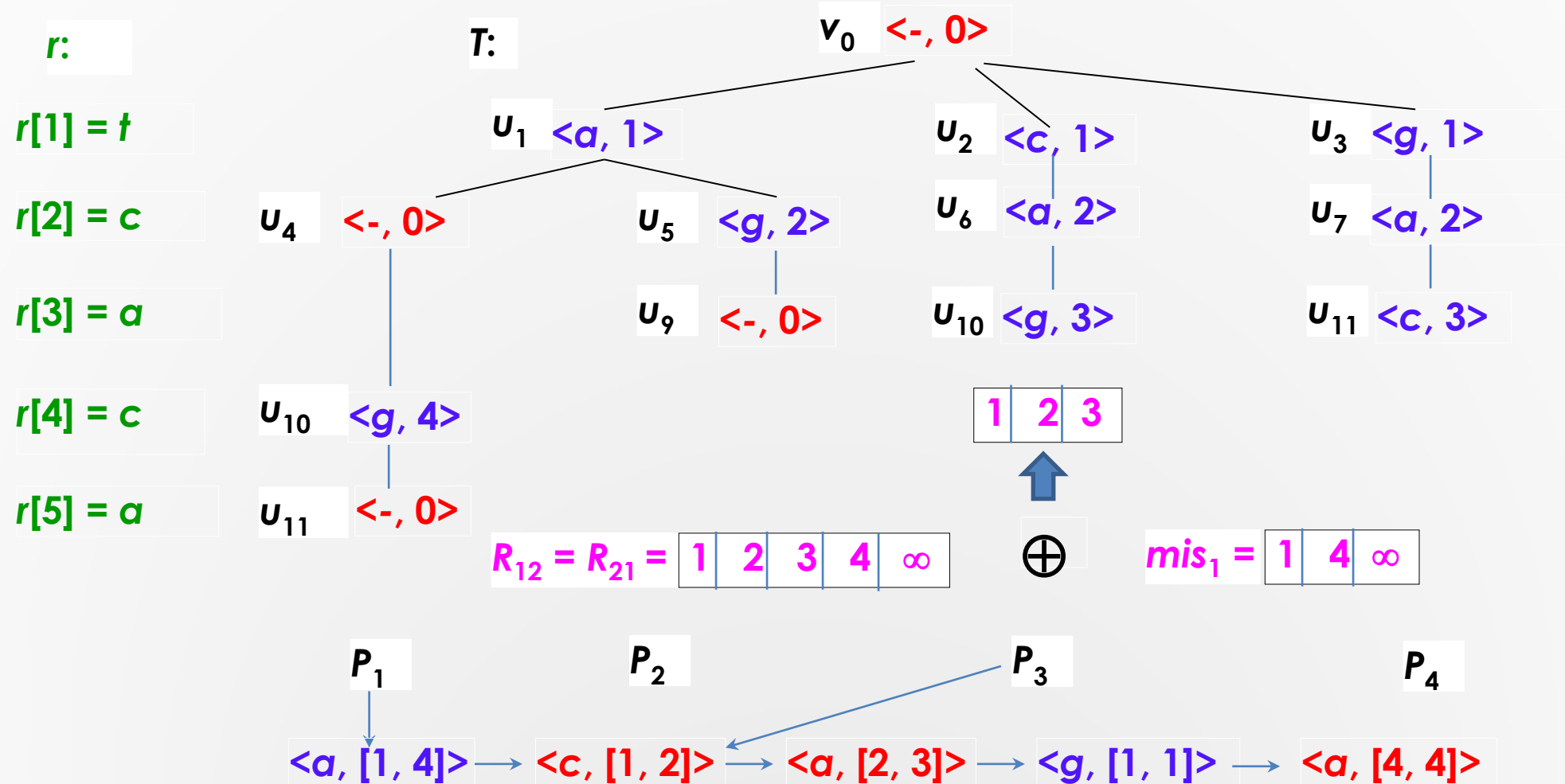


String Matching with k Mismatches

➤ Algorithm

- ❑ Mismatching tree generation
- ❑ Derivation of mismatching information for paths

Derivation of Mismatching Information



Algorithm

➤ Generation of mismatching trees

- In order to generate a mismatching tree D , we will use a stack S to control the process, in which each entry is a quadruple (v, j, κ, u) , where
 - v – a node inserted into the hash table.
 - j – j is an integer to indicate that v is the j th node on a path in T (counted from the root with the root as the 0th node).
 - κ – the number of mismatches between the path and $r[0 .. j]$ (recall that $r[0] = '-'$).
 - u – the parent of a node in D to be created for v .

Algorithm

❑ Mismatching tree generation

- ❑ Each time an entry $e = (v, j, \kappa, u)$ with $v = \langle x, [\alpha, \beta] \rangle$ is popped out from S , we will check whether $x = r[j]$.
 - If $x = r[j]$, we will generate a node $u' = \langle x, j \rangle$ and link it to u as a child.
 - If $x \neq r[j]$, we will check whether u is a node of the form $\langle -, 0 \rangle$. If it is not the case, generate a node $u' = \langle -, 0 \rangle$.
 - Otherwise, set u' to be u .
 - Using $\text{search}()$ to find all the children of v : v_1, \dots, v_l . Then, push each $(v_i, j + 1, \kappa', u')$ into S with κ' being κ or $\kappa + 1$, depending on whether $y_i = r[j + 1]$, where $v_i = \langle y_i, [\alpha_i, \beta_i] \rangle$.

Algorithm

❑ Mismatching information derivation for paths

- ❑ As with the basic process, each time a node $v = \langle x, [\alpha, \beta] \rangle$ (compared to $r[j]$) is encountered, which is the same as a previous one $v' = \langle x', [\alpha', \beta'] \rangle$ (compared to $r[i]$), we will not create a subtree in T in a way as for v' , but create a new node u for v in D (mismatching tree) and then go along $p(v')$ (the link associated with v') to find the corresponding nodes u' in D and search $D[u']$ in the breadth-first manner to generate a subtree rooted at u in D by simulating the merge operation discussed in Subsection B .
- ❑ In other words, $D[u]$ (to be created) corresponds to the mismatch arrays for all the paths going through v in T , which will not be actually produced.

Algorithm

❑ Mismatching information derivation for paths

- ❑ To this end, a queue data structure Q is used to do a breadth-first search of $D[u']$, and at the same time generate $D[u]$. In Q , each entry e is a pair (w, γ) with w being a node in $D[u']$, and γ an entry in R_{ij} . Initially, put $(u', R_{ij}[1])$ into Q , where $u' = \langle x, i \rangle$. In the process, when e is dequeued from Q (taken out from the front), we will make the following operations (simulating the steps in *merge()*):
 - 1) Let $e = (w, R_{ij}[l])$. Assume that $w = \langle z, f \rangle$ and $R_{ij}[l] = val$. If $\langle z, f \rangle = \langle -, 0 \rangle$, then create a copy of w added to $D[u]$. If w is not a leaf node, let w_1, \dots, w_h be the children of w and enqueue $(w_1, R_{ij}[l]), \dots, (w_h, R_{ij}[l])$ into Q (append at the end) in turn. If $\langle z, f \rangle \neq \langle -, 0 \rangle$, do (2), (3), or (4).
 - 2) If $f < i + val - 1$, add $\langle z, f - i + j \rangle$ to $D[u]$. If w is not a leaf node, enqueue $(w_1, R_{ij}[l]), \dots, (w_h, R_{ij}[l])$ into Q .

Algorithm

❑ Mismatching information derivation for paths

- 3) If $f = i + val - 1$, we will distinguish between two subcases: $z \neq r[j + val - 1]$ and $z = r[j + val - 1]$. If $z \neq r[j + val - 1]$, we have a mismatching and add a node $\langle z, j + val - 1 \rangle$ to $D[u]$. If $z = r[j + val - 1]$, create a node $\langle -, 0 \rangle$ added to $D[u]$. (If its parent is $\langle -, 0 \rangle$, it should be merged into its parent.)
- 4) If w is not a leaf node, enqueue $\langle w_1, R_{ij}[l + 1] \rangle, \dots, \langle w_h, R_{ij}[l + 1] \rangle$ into Q .
- 5) If $f > i + val - 1$, we will scan R_{ij} starting from $R_{ij}[l]$ until we meet $R_{ij}[l']$ such that $f \leq i + R_{ij}[l'] - 1$. For each $R_{ij}[g]$ ($l \leq g < l'$), we create a new node $\langle r[j + R_{ij}[g] - 1], j + R_{ij}[g] - 1 \rangle$ added to $D[u]$. Enqueue $\langle w, R_{ij}[l'] \rangle$ into Q .

Experiments

- **Compare 4 different approaches**
 - ❑ *BWT-based* [34] (BWT for short),
 - ❑ *Amir's method* [2] (Amir for short),
 - ❑ *Cole's method* [14] (Cole for short),
 - ❑ *Algorithm A discussed in this paper* ($A()$ for short)

Experiments

Test Environments:

- Implementation in C++, compiled by GNU make utility with optimization of level 2
- 64-bit Ubuntu operating system
- run on a single core of a 2.40GHz Intel Xeon E5-2630 processor with 32GB RAM

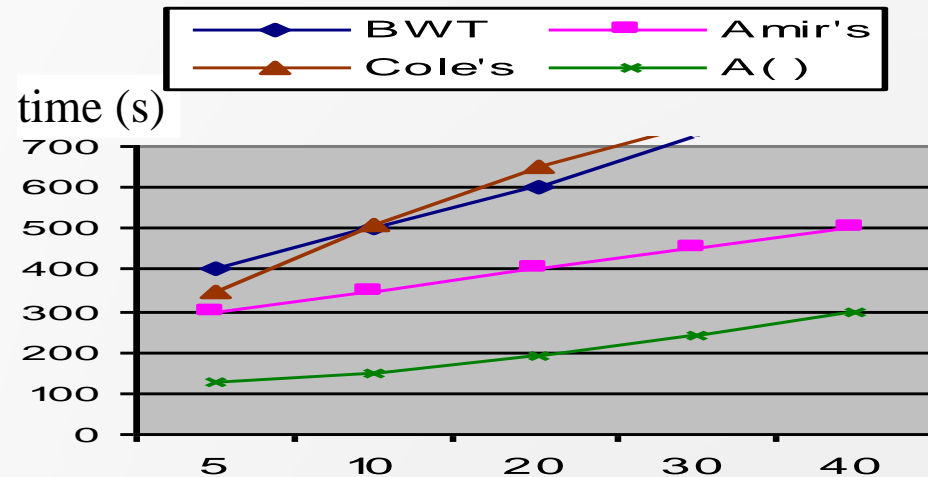
Experiments

TABLE I. CHARACTERISTICS OF GENOMES

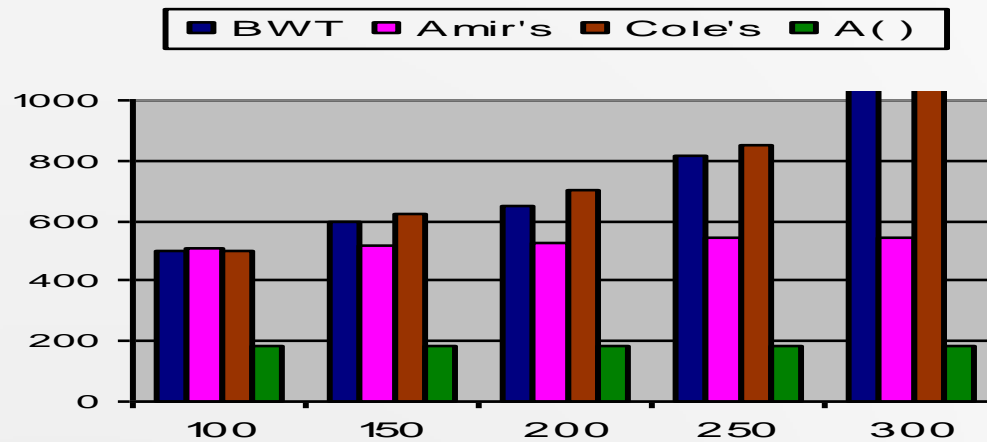
Genomes	Genome sizes (bp)
Rat chr1 (Rnor_6.0)	290,094,217
<i>C. merolae</i> (ASM9120v1)	16,728,967
<i>C. elegans</i> (WBcel235)	103,022,290
Zebra fish (GRCz10)	1,464,443,456
Rat (Rnor_6.0)	2,909,701,677

Tests with Real Data

➤ TESTS WITH VARYING LENGTH OF READS (OVER Rat genome)



varying values of k



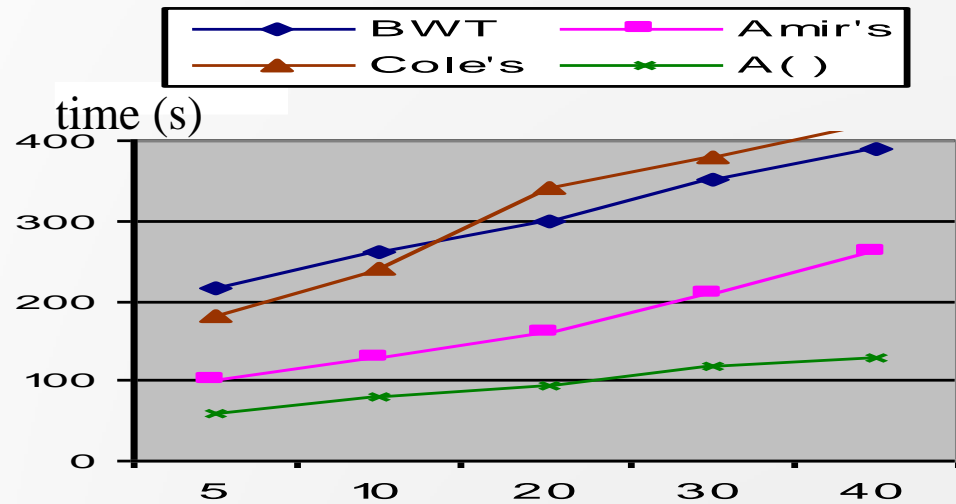
varying length of reads

Number of leaf nodes of S -trees

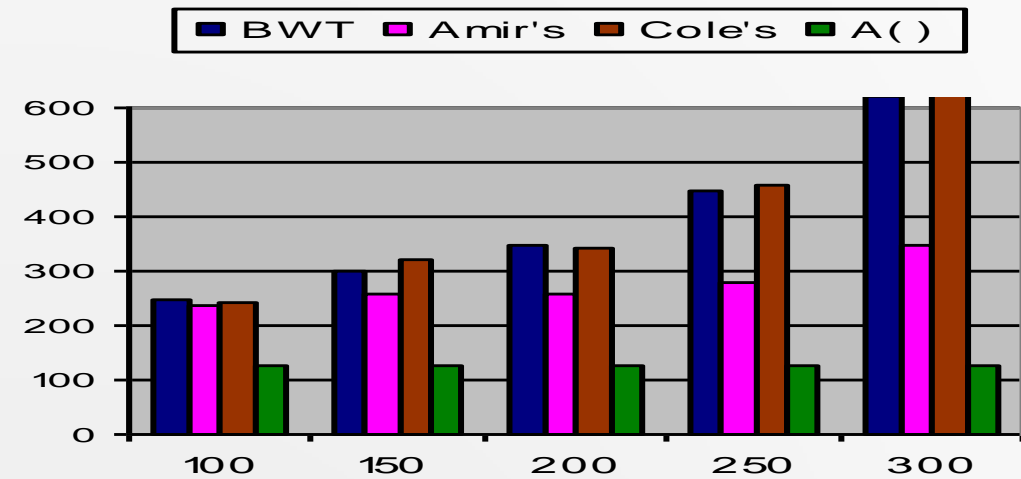
k/length-of-read	5/50	10/100	20/150	30/200
No. of leaf nodes	2K	0.7M	16.5M	102M

Tests with Real Data

➤ TESTS WITH VARYING LENGTH OF READS (OVER Zebra fish)



varying values of k



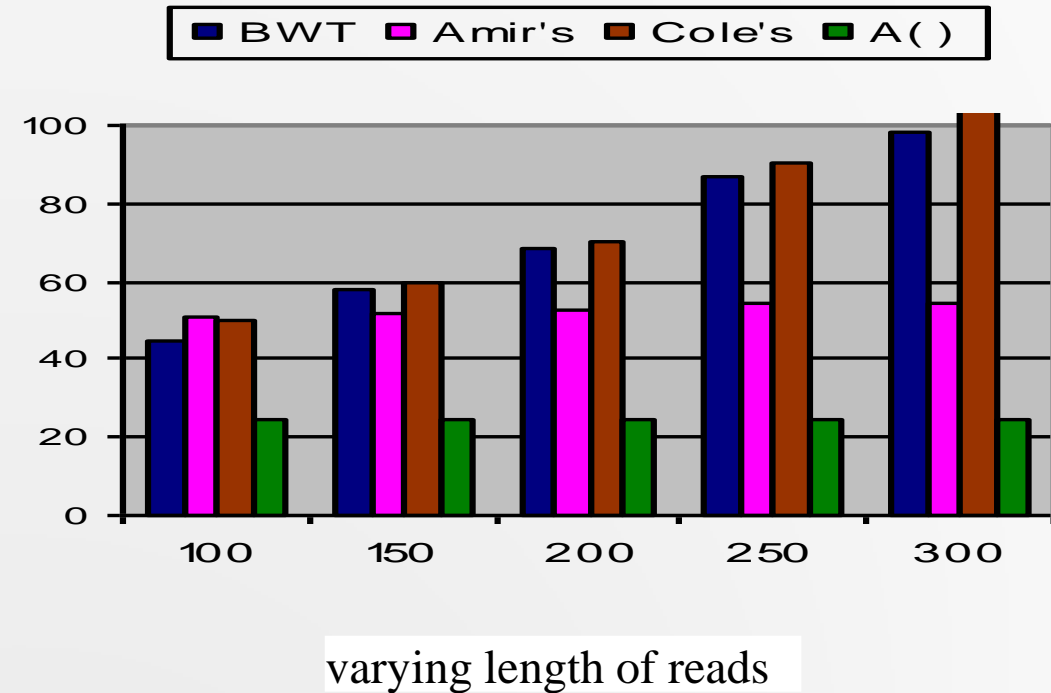
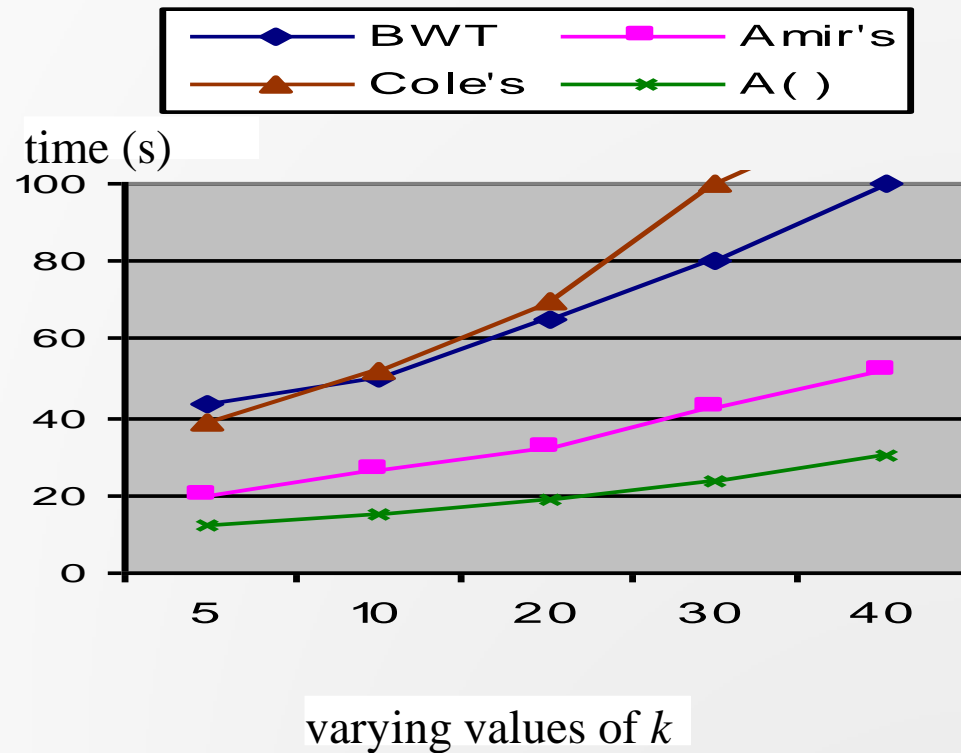
varying length of reads

Number of leaf nodes of S -trees

k/length-of-read	5/50	10/100	20/150	30/200
No. of leaf nodes	0.7K	0.30M	9.2M	89M

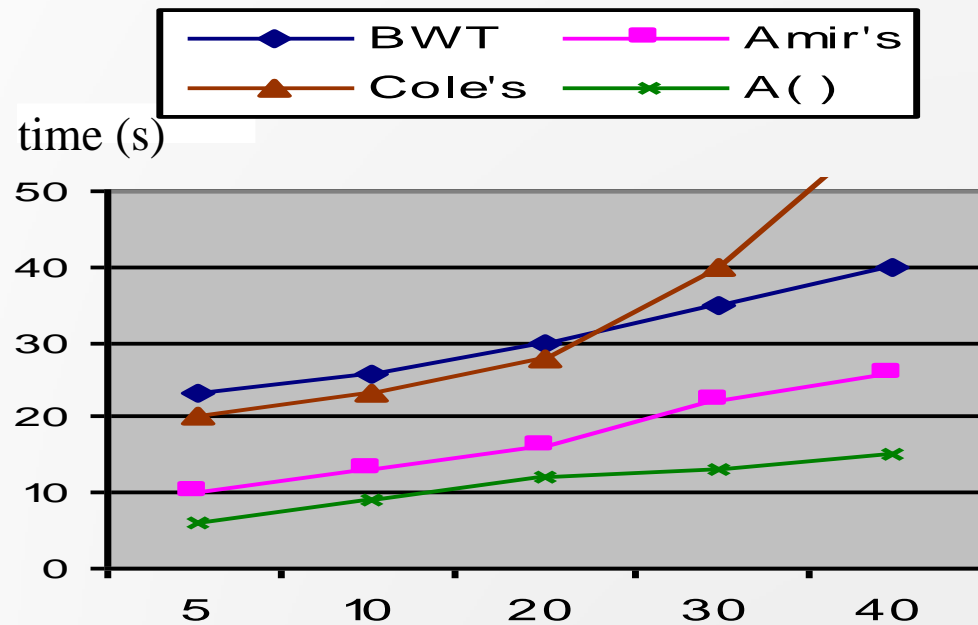
Tests with real Data

➤ Tests with varying length of reads (OVER Rat chr1)

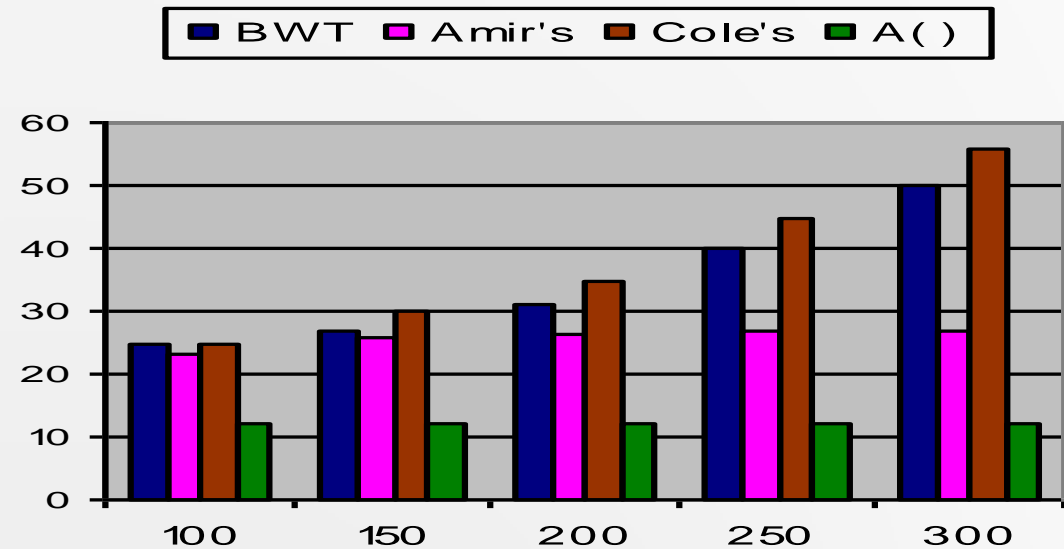


Tests with Real Data

➤ Tests with varying length of reads (OVER *C. elegans*)



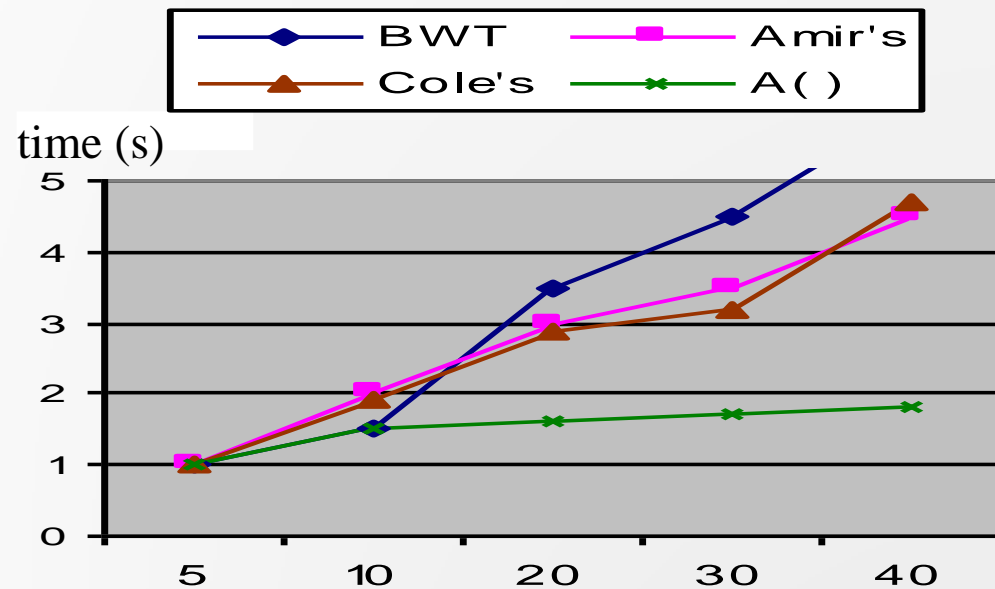
varying values of k



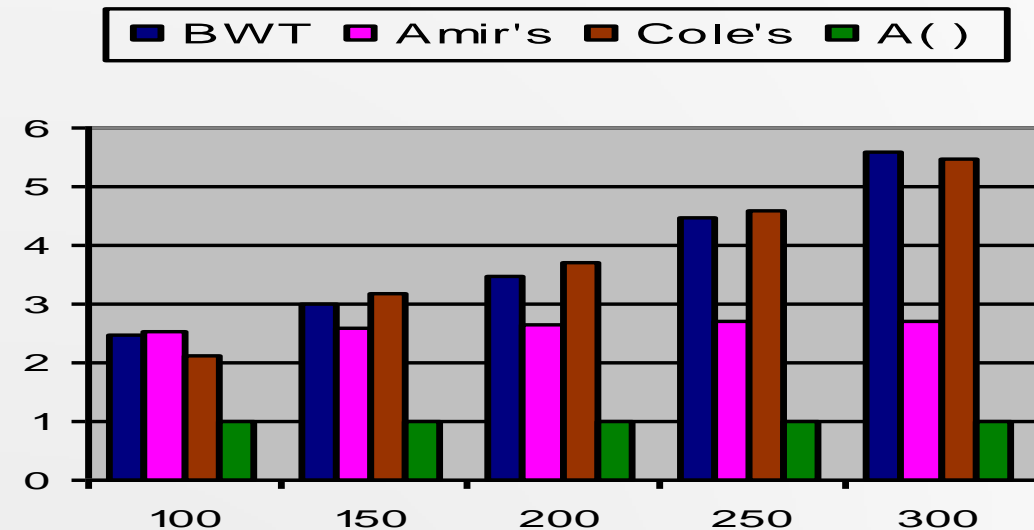
varying length of reads

Tests with Real Data

➤ Tests with varying length of reads (over *C. merlae*)



varying values of k



varying length of reads

Conclusion and Future Work

➤ Main contribution

- Combination of derivation of mismatching information and BWT indexes for k mismatching problem
- Concept of mismatching trees
- Extensive tests

➤ Future work

- String matching with k differences
- String matching with *don't care* symbols

Thank you!