

# התאמת מחרוזות:

- נתונים טקסט  $T$  ותבנית  $P$  שניהם הם אלפבית  $\Sigma$
- מצא הוק  $m$  היא זר מחרוזת של  $T$  (find)

- במכה אחר:  $O(|T| + |P|)$  [Knuth Morris Pratt]

- מחרוזת נתונת סטרי (Indexing):

- נבנה פעם אחר ואז נשאל על היכה תבנית

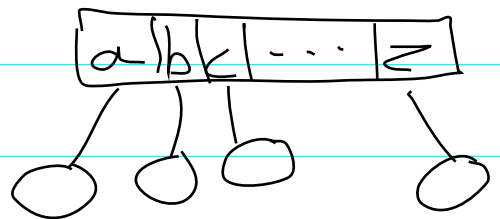
מקום  $O(|T|)$  הוק  $O(|T|)$

זמן בנייה  $O(|T| + \sum_{i \in \Sigma} f_i)$   $O(|T|)$  בזמן  $\Sigma$  זקום אנוליר

זמן שאילוח  $O(|P|)$  DNA/

## Trie מכנה לטמיה מילון

- לדל צורה  $\Sigma$  ילקים (לדל היא  $\Sigma$ )

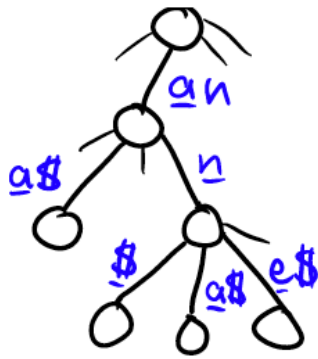
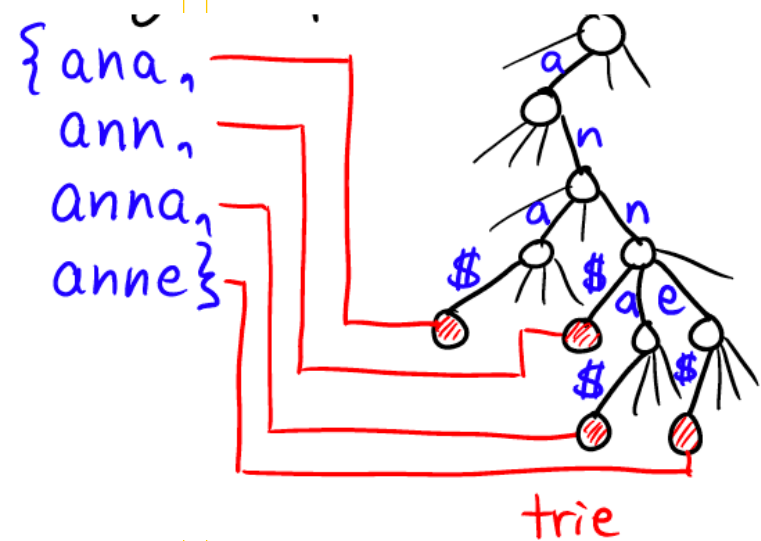


מסלול שורה - עדה

מילון מילון

- נסיים על מילון: בדו

אקט \$ להבדיל בין מילים לר"סור

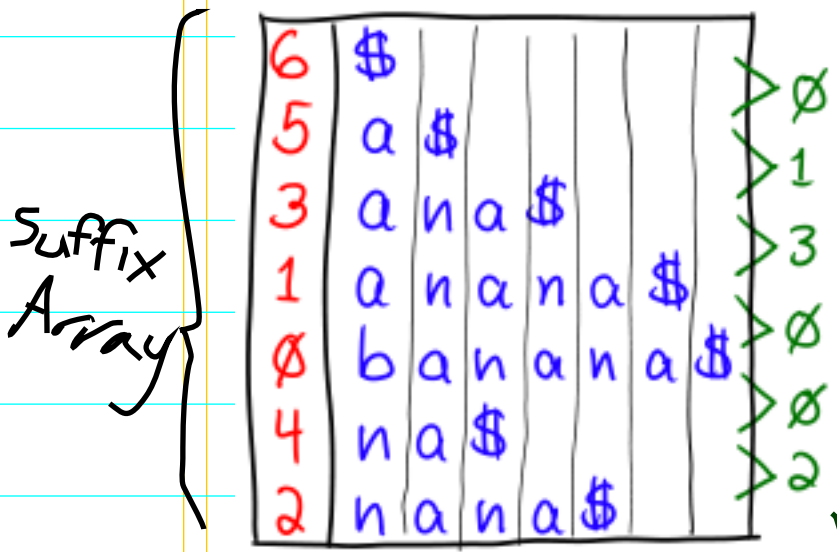


compressed trie

הכנסת מילים "מכונן" - Compressed Trie  
 המבנה הזה הוא מהיר וקטן, ויש לו יתרון

## Suffix Tree

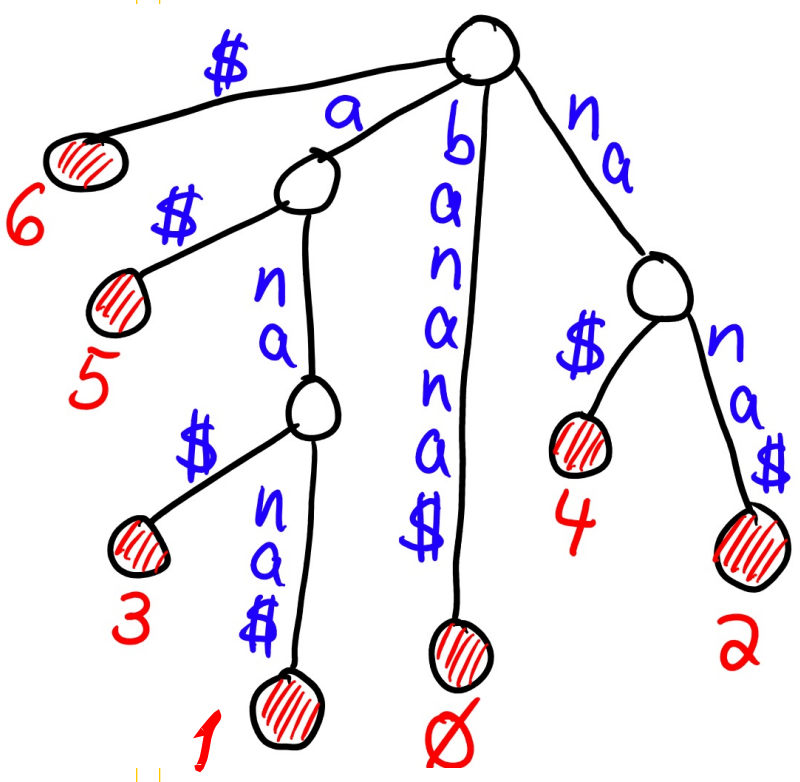
T\$ הוא מיון |T|+1 של ה-Compressed Trie



banana\$ : 0123456  
 0 1 2 3 4 5 6

- מיון של |T|+1
- מיון של ה-Compressed Trie
- T-סדר המיון של "א"

### LCP Array



- זמן מיון  $O(|T|)$
- זמן חישוב  $O(|T|)$
- זמן חישוב  $O(|T|)$

# Suffix Array: האינדקס של הסיומת בארזי מיון

- העלות של ה-suffix-tree ממשלה לימין
- ניתן ליישם בו בקנה אחד עם  $O(|P| \log |T|)$  במסמך
- ס"י מיון בינארי

• באי יתק אטק מע"ק ה-LCP (כיבודק)  
 כמה חוליק שהיק יש בתחילת  
 הסיומת ה-i וה-j בקשר המיון

$$LCP(i, j) = \text{RMQ}(LCP[i], \dots, LCP[j])$$

↑  
range minimum query

• בעזרת SA, LCP, RMQ ניתן לתבוע  
 כל פ עם מסמך  $O(|P| + \log |T|)$

- נניח כי מתקיים  $P = bacde$  במיון בינארי  
 SA של  $P$ , ונניח כי  $k=2$  חוליק

	$1 = LCP(L, M) < k$				$3 = LCP(L, M) > k$		
SA:	L	M	R		L	M	R
$k=2$	b	b			b	b	
	a	c			a	a	
	a	e			a	a	
	\$	\$			\$	c	
		\$				\$	

$LCP(L, M) = k$  - אורך המשותף הארוך ביותר  
 בין שתי מחרוזות  $L$  ו- $M$   
 -  $O(n \log n)$  - מורכב מ-  
 -  $O(n)$  - מורכב מ-

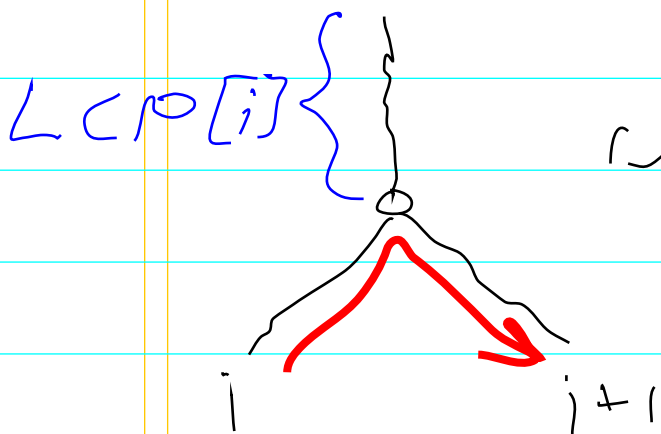
## סקירת SA ו-ST

$ST \leq SA + LCP$ : פשוט לראות את העלים ב-ST  
 משולש עינין

ה-LCP הוא עומק ה-LCA של שני ענפים עוקבים  
 באינדיקס  $\rightarrow$  באינדיקס  $\rightarrow$  קבוע יחיד

$ST \leq SA + LCP$ : מכניסים את היסודות ל-ST

אחת אחת עם סדר ב-SA  
 - ע"י עינין עוקבים. משלבים ב-ST מהעלה  
 פנויה עם LCP



• לא משלבים את האותיות  
 צומת פעמיים

$\Downarrow$   
 $O(n)$

$O(|T| + \text{Sort}(\Sigma))$  ממש (LCP+ SA) א"נ

[Kärkkäinen & Sanders 2003]

$T = \text{abcbbadeabca}$

$O(|T|)$  Radix Sort ממש  $\rightarrow$   $\text{Sort}(\Sigma)$  ממש  $\Sigma$  א"נ (1)

$T_1 = \underbrace{232}_{232} \underbrace{214}_{214} \underbrace{511}_{511} \underbrace{31}_{31}$   
 $T = 123221451231$   
 $T_0 = \underbrace{123}_{123} \underbrace{221}_{221} \underbrace{451}_{451} \underbrace{231}_{231}$

(2) מתייחס רקורסיבית עם ה"ט"  $T_0, T_1$  ו"נ"  $n$  "איות" (מיון האותיות עם "Radix")

$\Leftarrow$  מקבלים את ה"ט" ממיון  $LCP +$  מ"ט

$$\text{Suffixes}(T) = \text{Suffixes}(T_0) + \text{Suffixes}(T_1) + \text{Suffixes}(T_2)$$

(3) מיון א"נ  $\text{Suffixes}(T_2)$  עם "Radix" אחר מנימוק כ"ס

מ"ט  $\text{Suffixes}(T_2)$  :  $\langle a, T_0 \text{ - ב"א} \rangle$

מיון מנימוק ה"ט"  $T_0 =$  מספר בין 1 ל- $\frac{n}{3}$  א"נ האות

• אחר המיון, ה- $LCP$  בין שתי מ"ט א"נ מנימוק ה"ט" :  
 $LCP$  בין שתי מ"ט א"נ  $T_0$  + מנימוק ה"ט"

$\text{Suffixes}(T_2)$  - Merge הטו SA אר סגסג (4)  
 :  $\text{Suffixes}(T_0) + \text{Suffixes}(T_1)$  פס

אפם זרז דהטוולת  $T_0$ -נ כא"ס פם  $T_2$ -נ כא"ס  
 $\langle a, T_0\text{-נ כא"ס} \rangle \quad \langle a, T_1\text{-נ כא"ס} \rangle$

אפם זרז דהטוולת  $T_1$ -נ כא"ס פס  $T_2$ -נ כא"ס  
 $\langle ab, T_1\text{-נ כא"ס} \rangle \quad \langle ab, T_0\text{-נ כא"ס} \rangle$

- באופן קוטה, נקבלים את הערך  $LCP$  ס"י אית או שר"ס  
 $T_2$  דטו  $T_0, T_1$  דטו  $LCP$  סמוס

$$T(n) = T\left(\frac{n}{2}\right) + O(n) = O(n) : n > 0$$