

הוכן ע"י אמיר רובינשטיין

# מבני נתונים ומבוא לאלגוריתמים

---

נושא 12

מחרוזות  
Strings

# בתוכנית

- נכיר מבנה נתונים יעיל למילון של מחרוזות - trie.
- נכיר אלגוריתם למציאת תת-סדרה משותפת ארוכה ביותר.  
זהו אלגוריתם הפועל בשיטת תכנון דינאמי (dynamic programming).  
(ביבליוגרפיה לנושא זה: מבוא לאלגוריתמים, פרק 16 במהדורה הראשונה)

## מבנה הנתונים trie

# מוטיבציה

הגדרה: מחרוזת (string) מעל קבוצה נתונה  $\Sigma$  היא סדרה של איברים מתוך  $\Sigma$ .

את  $\Sigma$  נכנה א"ב, ואת איבריה אותיות.

על מחרוזות אפשר להגדיר יחס סדר, שנקרא סדר לקסיקוגרפי (lexicographic order):

נאמר שהמחרוזת  $a = a_0a_1\dots a_p$  קטנה לקסיקוגרפית מהמחרוזת  $b = b_0b_1\dots b_q$

אם מתקיים אחד משני התנאים:

1. קיים שלם  $j$  כך ש-  $a_i = b_i$  עבור  $i=0,1,\dots,j-1$  וגם  $a_j < b_j$ .

או:

2.  $p < q$  ו-  $a_i = b_i$  עבור  $i=0,1,\dots,p$ .

# מוטיבציה

המשימה - לממש ביעילות מילון (כולל מינימום, עוקב, וכו') של מחרוזות.

אפשרויות מוכרות: AVL - (שבכל צומת שלו יש מחרוזת או מצביע למחרוזת)

- hash table (שבכל איבר שלה יש מחרוזת או מצביע למחרוזת)

אפשר למשל להשתמש בטבלת גיבוב עם שרשור. נניח שבידנו פונקציה גיבוב טובה למחרוזות.

נתחו את יעילותם של המימושים הללו.

זיכרו: השוואה בין שתי מחרוזות דורשת זמן ליניארי באורך המחרוזת הקצרה מהשתיים.

- הניחו לשם פשטות שכל המחרוזות באותו אורך  $l$ .

- מספר המחרוזות יסומן כרגיל ב-  $n$ .

# מוטיבציה

## AVL

- הכנסה / חיפוש: בכל צומת תתבצע השוואה בזמן  $\Theta(l)$ , ומספר ההשוואות הנ"ל הוא  $\Theta(\log n)$ .  
לכן כל הפעולות יבוצעו בזמן  $\Theta(l \cdot \log n)$ .
- מחיקה / מינימום / עוקב:  $\Theta(\log n)$ .

## hash table

- הכנסה:  $\Theta(l)$  במקרה הגרוע, בהנחה שחישוב ערך הגיבוב למחרוזת דורש  $\Theta(l)$  (אחרת  $\Theta(1)$ ).
- מחיקה: בזמן קבוע במקרה הגרוע, אם הרשימות דו-כיווניות ונתון מצביע.
- חיפוש:  $\Theta(l)$  בממוצע,  $\Theta(l \cdot n)$  במקרה הגרוע.
- מינימום / עוקב:  $\Theta(l \cdot n)$  בכל המקרים

## מוטיבציה

נראה כעת מימוש יעיל יותר, שבו סיבוכיות הפעולות כלל לא תלויה במספר המחרוזות שבמבנה, אלא רק באורך.

# מבנה הנתונים trie

מקור השם trie הוא המילה retrieval (אחזור).

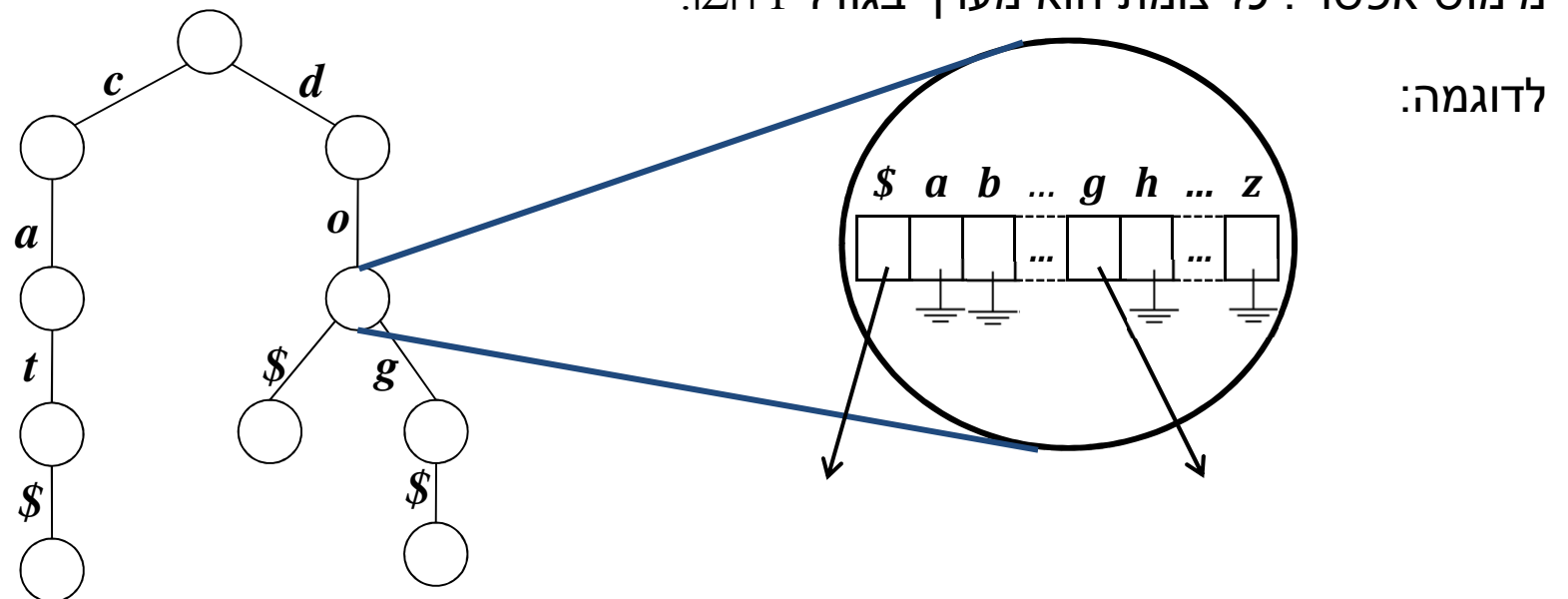
שתי הנחות:

1. לכל אורך ההרצאה נניח כי גודלו של  $\square$  קבוע לצורכי ניתוח סיבוכיות.

2. קיים תו \$ אינו שייך לא"ב  $\square$ . תו זה ישמש כקטן ביותר לקסיקוגרפית, וייצג סוף מחרוזת.

trie הוא עץ, שבו לכל צומת לכל היותר  $|\Sigma|+1$  בנים, עבור כל אחת מאותיות הא"ב ו-\$.

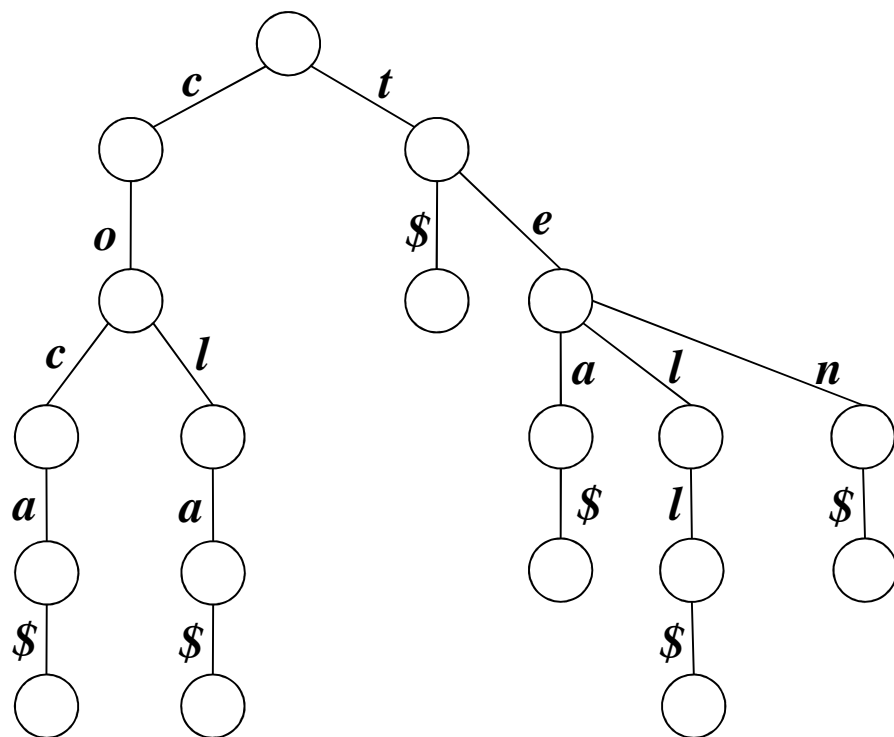
מימוש אפשרי: כל צומת הוא מערך בגודל  $|\Sigma|+1$ .





# מבנה הנתונים trie

עוד דוגמה:



אילו מחרוזות נמצאות ב- trie הבא?

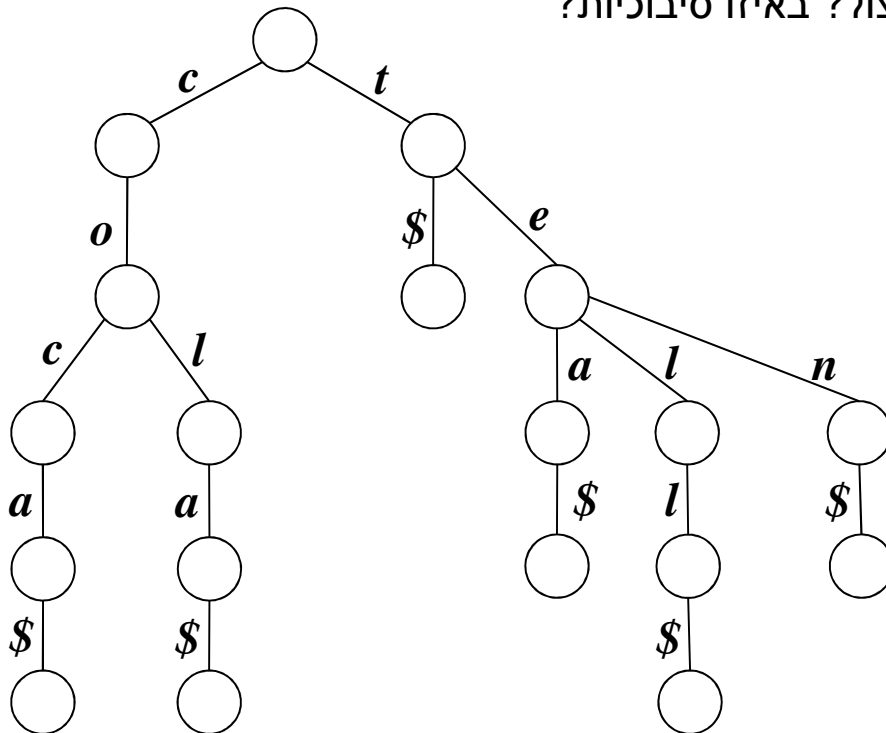
כל מחרוזת מיוצגת ע"י מסלול מהשורש לעלה (מספר העלים הוא כמספר המחרוזות).

נתאר כעת כיצד יבוצעו פעולות המילון השונות על trie.

# פעולות על trie

- Search( $s$ ) – עוקבים אחר המסלול של  $s$  ב- $s$ . trie נמצאת אם"ם אפשר להגיע לעלה ( $\$$ ).
- Insert( $s$ ) – עוקבים אחר המסלול של  $s$ , עד שנתקלים ב- $\text{Nil}$ , ואז ביתרת המסלול מקצים צמתים חדשים ומצביעים בהתאם. בסוף מוסיפים  $\$$ .
- Delete( $s$ ) – יש למחוק את חלק המסלול של  $s$  שאחרי נקודת הפיצול\* האחרונה שלו.

\* איך בודקים אם צומת הוא פיצול? באיזו סיבוכיות?



הפעולות הנ"ל רצות בזמן  $\Theta(|s|)$ .

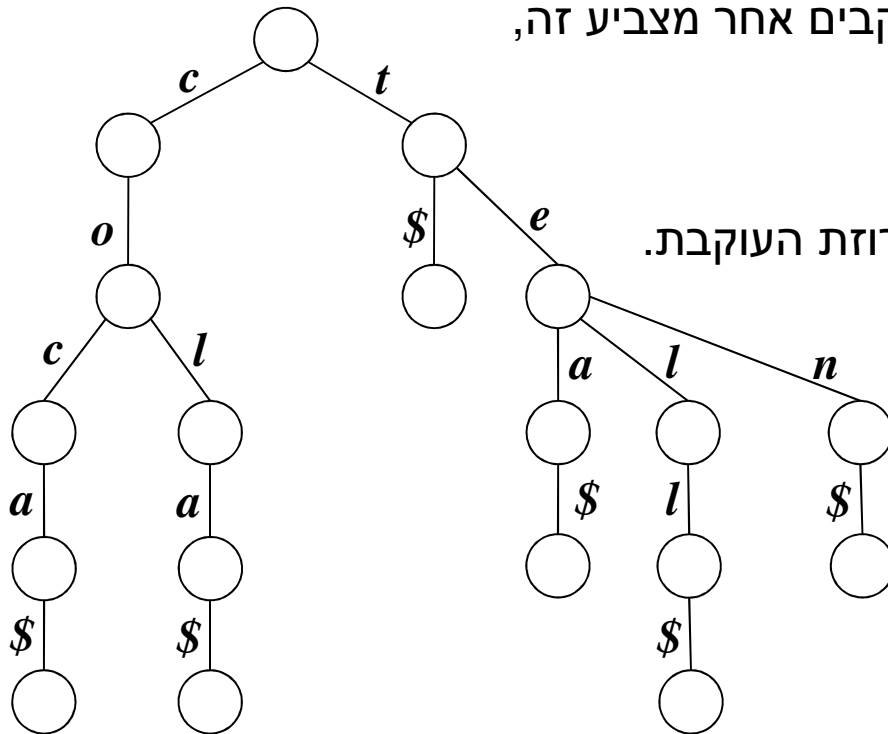
# פעולות על trie

- Minimum() – עוקבים אחר המסלול השמאלי ביותר, ו"רושמים את האותיות דרך עברנו".

סיבוכיות: בכל צומת יש לסרוק את המערך כדי לגלות מיהו המצביע השמאלי ביותר שאינו Nil.  
 פעולה זו דורשת  $\Theta(|\Sigma|) = \Theta(1)$  בכל צומת במסלול.  
 סה"כ  $\Theta(|s_{min}|)$  כאשר  $s_{min}$  היא המחרוזת המינימלית.

- Successor( $s$ ) – עוקבים אחר המסלול של  $s$  עד לעלה. משם חוזרים לנקודת הפיצול האחרונה שיש בה מצביע ימני יותר שאינו Nil, עוקבים אחר מצביע זה, ומשם ממשיכים כמו ב-Minimum.

סיבוכיות:  $\Theta(|s| + |s'|)$  כאשר  $s'$  היא המחרוזת העוקבת.



## סיבוכיות מקום

נתון trie המכיל  $n$  מחרוזות שאורכן הכולל  $m$ .

שאלה: מהי סיבוכיות הזיכרון הדרושה?

תשובה: מספר הצמתים כפול כמות הזיכרון שדורש כל צומת.

- כל צומת מכיל מערך בגודל  $|\Sigma|+1 = \Theta(1)$ .
- כמות הצמתים: ישנו שורש אחד, ובנוסף מחרוזת  $s_i$  מיוצגת ע"י  $|s_i|+1$  צמתים. במקרה הגרוע, אם אין צמתים משותפים כלל, כמות הצמתים הכוללת היא:

$$\# \text{צמתים} = 1 + \sum_{i=1}^n (|s_i| + 1) = 1 + m + n = \Theta(m)$$

סיבוכיות הזיכרון הדרושה היא אם כן ליניארית באורך הכולל  $m$ , ולא תלויה כלל במספר המחרוזות.

## מיון מחרוזות

נתונות  $n$  מחרוזות, שאורכן הכולל הוא  $m$ .  
אנו מחפשים אלגוריתם יעיל שממין את המחרוזות לפי סדר לקסיקוגרפי.

פתרון אפשרי (לא יעיל):

להכניס את המחרוזות ל-AVL (בכל צומת מחרוזת או מצביע למחרוזת), ולסייר בעץ inorder.

סיבוכיות:

נניח לשם פשטות שכל מחרוזת באורך אחיד  $m/n$ .

הכנסת  $n$  המחרוזות לעץ דורשת  $\Theta\left(\frac{m}{n} \sum_{i=1}^n \log i\right) = \Theta(m \log n)$   
הסיור כולל ההדפסות  $\Theta(n+m) = \Theta(m)$ .

סה"כ:  $\Theta(m \log n)$

נראה כעת כיצד שימוש ב-trie מוריד את הסיבוכיות ל- $\Theta(m)$ .

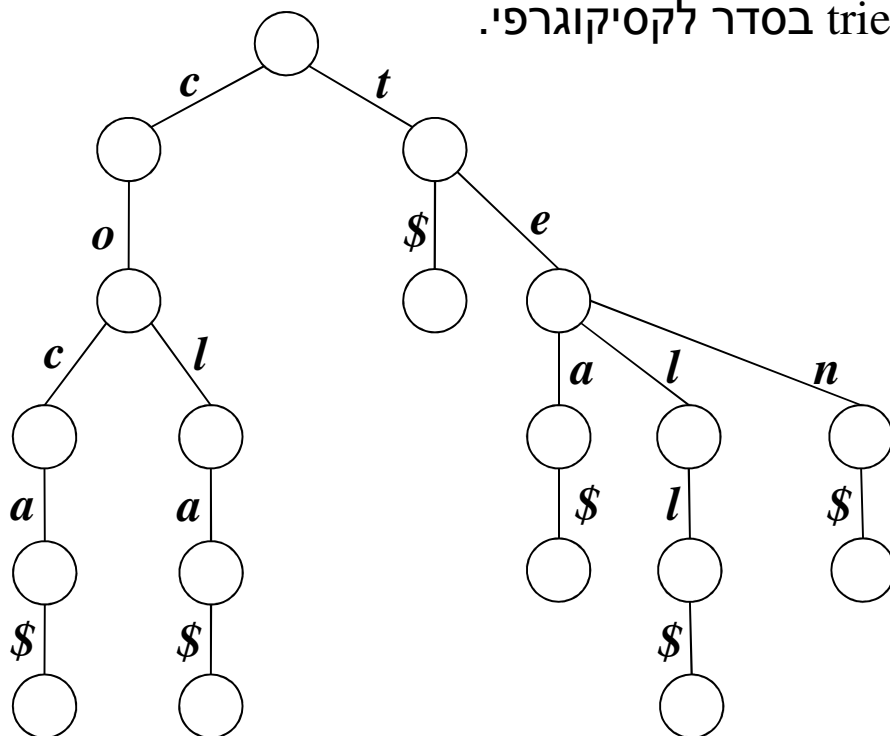
# מיון מחרוזות בעזרת trie

## תרגיל

נתון trie כלשהו, ובו  $n$  מחרוזות שאורכן הכולל הוא  $m$ .

תארו אלגוריתם שמדפיס את המחרוזות שב- trie בסדר לקסיקוגרפי.

מה סיבוכיות הפתרון שלכם?



# מיון מחרוזות בעזרת trie

בהמשך לתרגיל הקודם, להלן אלגוריתם למיון  $n$  מחרוזות, שאורכן הכולל הוא  $m$ .

1. מכניסים ל-trie את המחרוזות בזו אחר זו.
2. מבצעים סיור pre-order ב-trie באופן הבא:
  - 2.1 מאתחלים רשימה מקושרת ריקה
  - 2.2 בכל ירידה בעץ מוסיפים את האות "דרכה עברנו" לסוף הרשימה
  - 2.3 בכל עלייה בעץ מוציאים מהרשימה את האות שבסופה
  - 2.4 בכל פעם שמגיעים לעלה, מדפיסים את תוכן הרשימה מתחילתה עד סופה

## סיבוכיות

$$1. \Theta\left(\sum_{i=1}^n |s_i|\right) = \Theta(m)$$

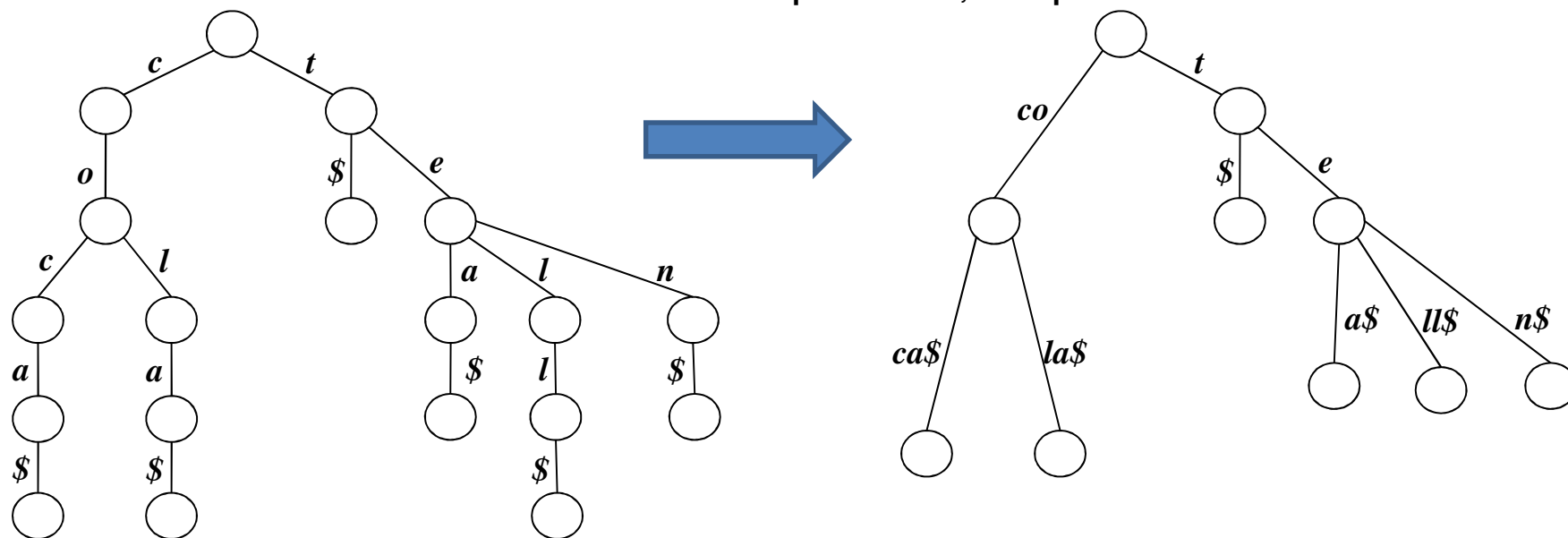
2. מספר הצמתים בעץ הוא כאמור במקרה הגרוע  $\Theta(m)$ , וסיור בעץ מתבצע כזכור בזמן ליניארי במספר הצמתים (שורות 2.2 ו-2.3 משפיעות רק על הקבועים).
- שורה 2.4 לא בהכרח רצה בזמן קבוע כל פעם, אבל סך כל זמן הריצה שלה הוא  $\Theta(m)$  (מדוע?)  
סה"כ  $\Theta(m)$ .

# שיפור סיבוכיות מקום - דחיסה

נתון trie המכיל  $n$  מחרוזות שאורכן הכולל  $m$ .

כזכור, סיבוכיות הזיכרון הדרושה היא ליניארית באורך הכולל  $m$ , ולא תלויה כלל במספר המחרוזות.  
האם ניתן לחסוך בזיכרון?

דחיסת trie - ביטול צמתים בעלי בן יחיד, כדי לחסוך בצמתים.



לאחר הדחיסה מספר הצמתים קטן מ-  $\Theta(m)$  ל-  $\Theta(n)$  :

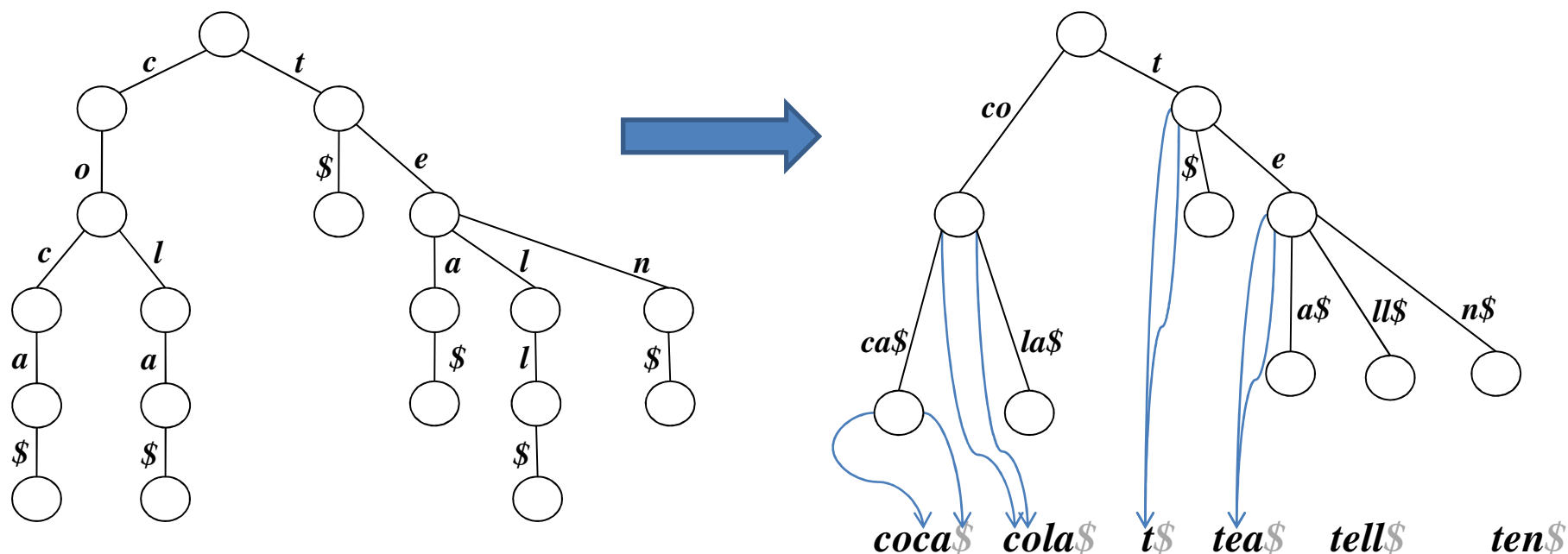
- יש  $n$  עלים (כמספר המחרוזות), לכל היותר  $n/2$  אבות, לכל היותר  $n/4$  סבים, וכו'.

- לכן מספר הצמתים לא גדול מ-  $n + n/2 + n/4 + \dots < 2n$ .



# שיפור סיבוכיות מקום - דחיסה

נוצרה בעיה – לא מספיק להחזיק מערך בכל צומת, כי צמתים מייצגים תתי-מחרוזות ולא אותיות.  
פתרון: בכל צומת נשמור שני מצביעים – לתחילת תת-המחרוזת שמייצג הצומת ולסופה.



כך, בכל צומת דרוש עדיין  $\Theta(1)$  זיכרון, וכאמור כמות הצמתים ירדה ל-  $\Theta(n)$ .  
אבל צריך כעת לשמור עותקים מהמחרוזות, אליהם יופנו המצביעים מהצמתים – סה"כ  $\Theta(m)$ .

מתי אם כן דחיסת trie היא בעלת ערך בחיסכון בזיכרון?

# שיפור סיבוכיות מקום בעץ סיומות

השיפור הנ"ל בסיבוכיות הזיכרון רלוונטי יותר בגרסה חשובה ושימושית של trie שנקראת "עץ סיומות" (suffix tree).

זהו trie אליו הוכנסו כל הסיומות של מחרוזת מסוימת.

לדוגמא, עץ סיומות עבור המחרוזת  $s = \text{"banana"}$ :

המחרוזות בעץ:

מחרוזת ריקה

a

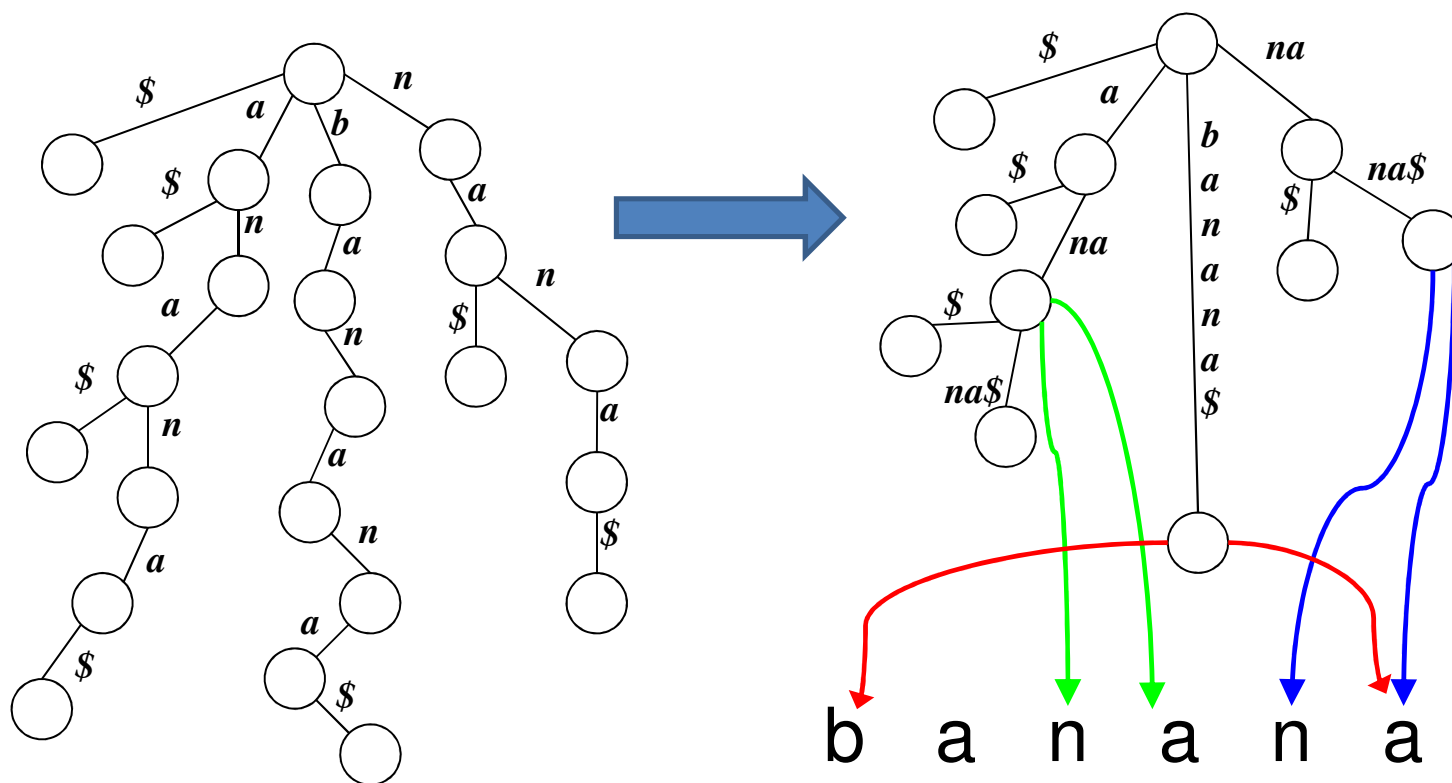
na

ana

nana

anana

$s = \text{bananas}$



## עץ סיומות

לעצי סיומות מספר שימושים חשובים, ביניהם:

1. מציאת תת-מחרוזת בתוך מחרוזת נתונה (מנוע חיפוש, ביואינפורמטיקה).
2. דחיסת אינפורמציה (למשל אלגוריתם ziv-lempel compression).
3. עבור שתי מחרוזות נתונות, מציאת תת-מחרוזת (רצופה) משותפת ארוכה ביותר.

ניתן לקרוא על כך למשל ב-

Algorithms on strings, Trees and sequences, Dan Gustfield

Chapter 5, 7.3, 7.4, 7.17

תת-סדרה משותפת ארוכה ביותר

Longest common subseries (LCS)

# תת-סדרה משותפת ארוכה ביותר

## הגדרת הבעיה

קלט: שתי מחרוזות:  $X$  באורך  $m$  ו-  $Y$  באורך  $n$ .  
פלט: אורך תת-סדרה משותפת ארוכה ביותר.

הערה: תת-סדרה משותפת לא חייבת להיות רצופה (בניגוד ל"תת-מחרוזת").

דוגמה:  $X = bacaacb$        $Y = acabc$

תתי סדרות משותפות מקסימליות:  $acab$  או  $acac$

פתרון נאיבי: לבדוק כל תת-סדרה אפשרית של  $X$ :  $\Theta(n2^m)$ . (כמה תת-סדרות יש לסדרה באורך  $m$ ?)

## סימונים

- תת-סדרה משותפת ארוכה ביותר של מחרוזות  $X$  ו-  $Y$  תסומן  $LCS(X, Y)$ .
- האורך של  $LCS(X, Y)$  יסומן  $L(X, Y)$ .
- עבור מחרוזת  $X$ :
  - נסמן ב-  $x_i$  את התו ה-  $i$  של  $X$
  - נסמן ב-  $X[i..j]$  את תת-הסדרה מ-  $x_i$  עד  $x_j$ .

## תת-סדרה משותפת ארוכה ביותר

3 אבחנות, שיעזרו לגיבוש פתרון

עבור שתי מחרוזות:  $X$  באורך  $m$  ו-  $Y$  באורך  $n$  מתקיים:

אבחנה 1: אם אחת הסדרות ריקה אז  $L(X,Y) = 0$  ו-  $LCS(X,Y) = \square$

אבחנה 2: אם  $x_m = y_n$  אז  $L(X,Y) = L(X[1..m-1], Y[1..n-1]) + 1$

$$LCS(X,Y) = LCS(X[1..m-1], Y[1..n-1]) \square x_m$$

אבחנה 3: אם  $x_m \neq y_n$  אז  $L(X,Y) = \max\{ L(X[1..m-1], Y) , L(X, Y[1..n-1]) \}$

$$LCS(X,Y) = LCS(X[1..m-1], Y) \quad OR \quad LCS(X, Y[1..n-1])$$

# פתרון רקורסיבי פשוט

לאור האבחנות הללו, להלן אלגוריתם רקורסיבי לחישוב  $L(X,Y)$

$L(X, Y)$

1.  $m \leftarrow \text{length}[X], n \leftarrow \text{length}[Y]$
2. **if**  $n=0$  or  $m=0$
3.       **return** 0
4. **if**  $x_m = y_n$
5.       **return**  $L(X[1..m-1], Y[1..n-1]) + 1$
6. **else**
7.       **return**  $\max \{ L(X[1..m-1], Y), L(X, Y[1..n-1]) \}$

## סיבוכיות

נסמן ב-  $s$  את סכום האורכים  $n+m$ .

במקרה הגרוע נכנסים לשורה 7 בכל פעם:

$$T(s) = 2T(s-1) + 1 = \Theta(2^s)$$

סיבוכיות אקספוננציאלית בסכום האורכים!

ממה זה נובע?

## פתרון משופר בשיטת "תכנון דינאמי"

בפתרונות רקורסיביים לבעיות, אם ישנן תת-בעיות חופפות זו לזו (כלומר יש להן תת-תת-בעיות משותפות), פעמים רבות הדבר גורר סיבוכיות זמן אקספוננציאלית. זאת מכיוון שהאלגוריתם פותר פעמים רבות את אותן תת-בעיות.

(בהקשר זה, היזכרו גם בפתרון הרקורסיבי לחישוב האיבר ה- $n$  בסדרת פיבונאצ'י).

תכנון דינאמי (dynamic programming) היא שיטה שמאפשרת חסכון בזמן ע"י כך שכל תת-בעיה נפתרת רק פעם אחת: מאחסנים בטבלה את הפתרונות שחושבו, ושולפים אותם מהטבלה בעת הצורך.



# הדגמת הפתרון בשיטת "תכנון דינאמי"

בטבלה הבאה,  $C[i, j]$  יכיל את  $L(X[1..i], Y[1..j])$ .

$Y =$	$a$	$c$	$a$	$b$	$c$
$X$	0	0	0	0	0
$  $					
$b$	0				
$a$	0				
$c$	0				
$a$	0				
$c$	0				
$b$	0				

אתחול: שורה ראשונה ועמודה ראשונה יכילו 0.

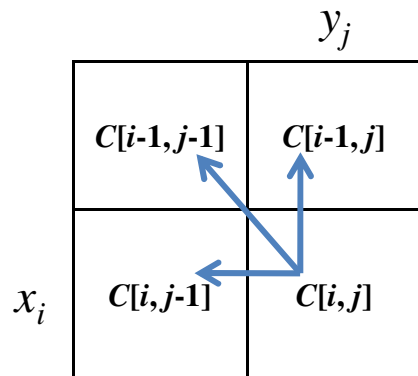
סיום: בסיום האלגוריתם הפלט יהיה רשום ב-  $C[m, n]$ .

כלומר  $L(X, Y) = C[m, n]$ .

צעדי האלגוריתם:

נחשב את הערכים בטבלה שורה אחר שורה.

אפשר לחשב כל תא מתוך 3 תאים סמוכים בלבד:



אם  $x_i = y_j$  אז  $C[i, j] = C[i-1, j-1] + 1$

אחרת -  $x_i \neq y_j$  אז

אם  $C[i-1, j] \geq C[i, j-1]$  אז  $C[i, j] = C[i-1, j]$

אחרת  $C[i, j] = C[i, j-1]$

# הדגמת הפתרון בשיטת "תכנון דינאמי"

אם  $x_i = y_j$  אז  $C[i, j] = C[i-1, j-1] + 1$

אחרת -  $x_i \neq y_j$  אז

אם  $C[i-1, j] \geq C[i, j-1]$  אז  $C[i, j] = C[i-1, j]$

אחרת  $C[i, j] = C[i, j-1]$

	$Y=$	$a$	$c$	$a$	$b$	$c$
$X$						
$\parallel$	0	0	0	0	0	0
$b$	0	$\uparrow 0$	$\uparrow 0$	$\uparrow 0$	$\swarrow 1$	$\leftarrow 1$
$a$	0	$\swarrow 1$	$\leftarrow 1$	$\swarrow 1$	$\uparrow 1$	$\uparrow 1$
$c$	0	$\uparrow 1$	$\swarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\swarrow 2$
$a$	0	$\swarrow 1$	$\uparrow 2$	$\swarrow 3$	$\leftarrow 3$	$\leftarrow 3$
$c$	0	$\uparrow 1$	$\swarrow 2$	$\uparrow 3$	$\uparrow 3$	$\swarrow 4$
$b$	0	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\swarrow 4$	$\uparrow 4$

## האלגוריתם

לשם נוחות נמספר שורות ועמודות החל ב-0.

$L(X, Y)$

```
1.  $m \leftarrow \text{length}[X], n \leftarrow \text{length}[Y]$ 
2. for  $i \leftarrow 0$  to  $m$        $C[i, 0] \leftarrow 0$       //leftmost column
3. for  $j \leftarrow 0$  to  $n$        $C[0, j] \leftarrow 0$       //top row
4. for  $i \leftarrow 1$  to  $m$ 
5.     for  $j \leftarrow 1$  to  $n$ 
6.         if  $x_i = y_j$ 
7.              $C[i, j] \leftarrow C[i-1, j-1] + 1$ 
8.         else if  $C[i-1, j] \geq C[i, j-1]$ 
9.              $C[i, j] \leftarrow C[i-1, j]$ 
10.        else       $C[i, j] \leftarrow C[i, j-1]$ 
11. return  $C[m, n]$ 
```

סיבוכיות  $\Theta(mn)$ .

שאלה: מה צריך להוסיף כדי שנוכל גם לשחזר את ה-LCS עצמה?

## שחזור תת-סדרה משותפת ארוכה ביותר

נשמור בטבלה נוספת  $B$  את ה"כיוונים".

בסיום האלגוריתם:

- נאתחל מחרוזת ריקה  $Z$
- נעקוב אחר המסלול החל ב-  $B[m, n]$ :
- בכל פעם ש-  $B[i, j] = \nwarrow$  נוסיף את  $x_i$  לתחילת המחרוזת  $Z$ .

שימו לב שבסיום  $Z = \text{LCS}(X, Y)$ , אבל ייתכן שיש עוד תת-סדרות ארוכות ביותר שונות מ-  $Z$ .

	$Y =$	<b><i>a</i></b>	<b><i>c</i></b>	<b><i>a</i></b>	<b><i>b</i></b>	<b><i>c</i></b>
$X$		0	0	0	0	0
$\parallel$		0	0	0	0	0
<b><i>b</i></b>		0	$\uparrow$ 0	$\uparrow$ 0	$\uparrow$ 0	$\nwarrow$ 1 $\leftarrow$ 1
<b><i>a</i></b>		0	$\nwarrow$ 1	$\leftarrow$ 1	$\nwarrow$ 1	$\uparrow$ 1 $\uparrow$ 1
<b><i>c</i></b>		0	$\uparrow$ 1	$\nwarrow$ 2	$\leftarrow$ 2	$\leftarrow$ 2 $\nwarrow$ 2
<b><i>a</i></b>		0	$\nwarrow$ 1	$\uparrow$ 2	$\nwarrow$ 3	$\leftarrow$ 3 $\leftarrow$ 3
<b><i>c</i></b>		0	$\uparrow$ 1	$\nwarrow$ 2	$\uparrow$ 3	$\uparrow$ 3 $\nwarrow$ 4
<b><i>b</i></b>		0	$\uparrow$ 1	$\uparrow$ 2	$\uparrow$ 3	$\nwarrow$ 4 $\uparrow$ 4

$$\text{LCS}(X, Y) = \text{"acac"}$$

סיבוכיות השחזור:  $\Theta(m+n)$

## שחזור תת-סדרה משותפת ארוכה ביותר

בנית הטבלה  $B$ :

$L(X, Y)$

```
1.  $m \leftarrow \text{length}[X], n \leftarrow \text{length}[Y]$ 
2. for  $i \leftarrow 0$  to  $m$        $C[i, 0] \leftarrow 0$       //leftmost column
3. for  $j \leftarrow 0$  to  $n$        $C[0, j] \leftarrow 0$       //top row
4. for  $i \leftarrow 1$  to  $m$ 
5.     for  $j \leftarrow 1$  to  $n$ 
6.         if  $x_i = y_j$ 
7.              $C[i, j] \leftarrow C[i-1, j-1] + 1$ 
8.         else if  $C[i-1, j] \geq C[i, j-1]$ 
9.              $C[i, j] \leftarrow C[i-1, j]$ 
10.        else       $C[i, j] \leftarrow C[i, j-1]$ 
11. return  $C[m, n]$ 
```

$B[i, j] \leftarrow \nwarrow$

$B[i, j] \leftarrow \uparrow$

$B[i, j] \leftarrow \leftarrow$

למעשה, היה אפשר לוותר על הטבלה  $B$ . כיצד?

## שאלות חזרה

$X = BABBDAC$

$Y = BBADCA$

1. מצאו תת-סדרה משותפת ארוכה ביותר של המחרוזות:

ע"י הרצת האלגוריתם LCS.

מיצאו גם תת-סדרות משותפות ארוכות ביותר אחרות מזו שהאלגוריתם מחזיר.

# תשובות לשאלות חזרה

1. *BADA, BADC, BBDA, BBDC, BBAC*

# תרגילים



תרגילים מומלצים מהספר

פרק 12 במהדורה השנייה (2008)  
בעיה 12-2

פרק 16 במהדורה הראשונה (1998)  
16.3-1  
16.3-2  
16.3-5

## תרגילים נוספים

1. הציעו אלגוריתם שזמן ריצתו  $O(n^2)$  אשר מוצא תת-סדרה לא יורדת ארוכה ביותר של סדרה בת  $n$  מספרים.

2. דרוש מבנה נתונים למימוש הפעולות:

Init( $\{s_1, s_2, \dots, s_n\}$ ) – אתחול המבנה בהינתן קבוצה של  $n$  מחרוזות.

סיבוכיות -  $O(m)$  כאשר  $m$  הוא האורך הכולל של כל המחרוזות.

Find-Permute( $s$ ) – בדיקה האם  $s$  היא פרמוטציה של אחת המחרוזות  $s_1, \dots, s_n$ .

סיבוכיות -  $O(|s|)$ .

כלומר, בהינתן קבוצת מחרוזות ידועה מראש, יש לענות על השאלה הבאה:  
האם מחרוזת חדשה כלשהי היא פרמוטציה של אחת המחרוזות הנתונות?

3. עץ סיומות (suffix tree) של מחרוזת נתונה  $s$  הוא trie אליו הוכנסו כל הסיומות של  $s$ .  
(סיומת של מחרוזת היא תת-מחרוזת שלה החל ממקום מסוים ועד סופה).

בהינתן עץ סיומות של מחרוזת  $s$ , הראו כיצד ניתן למצוא בעזרתו האם מחרוזת  $s'$  היא תת-מחרוזת (רצופה) של  $s$ .

## פתרון 1

עבור סדרה  $X$ , ניצור עותק שלה  $X'$  ונמין אותו. כעת נקרא ל-  $\text{LCS}(X, X')$ . התוצאה היא תת-סדרה ארוכה ביותר המשותפת לסדרה המקורית ולעותק הממוין שלה, וזוהי גם תת-סדרה מונוטונית ארוכה ביותר של הסדרה המקורית.

זמן ריצה:  $\Theta(n \log n)$  עבור המיון, ועוד  $\Theta(n^2)$  עבור האלגוריתם  $\text{LCS}$ . סה"כ  $\Theta(n^2)$ .

## פתרון 2

:Init

- נמין כל מחרוזת  $s_i$  באמצעות Counting-Sort.  
סיבוכיות מיון מחרוזת בודדת  $s_i$ :  $\Theta(n+k) = \Theta(|s_i| + |\Sigma|) = \Theta(|s_i|)$   
סה"כ  $\Theta(m)$
- נכניס את כל המחרוזות הממוינות ל- trie בזו אחר זו.  
סיבוכיות הכנסת מחרוזת בודדת  $\Theta(|s_i|)$ , ובסה"כ  $\Theta(m)$ .

:Find-Permute

- נמין את  $s$  באותו אופן, ונבצע חיפוש ב- trie של  $s$  הממוינת. סיבוכיות המיון + החיפוש  $\Theta(|s|)$ .

## פתרון 3

תת-מחרוזת של מחרוזת נתונה היא התחלה (רישא) של סיומת (סיפא) כלשהי שלה.  
לכן  $s'$  היא תת-מחרוזת של  $s$  אם"ם קיים מסלול של  $s'$  שמתחיל בשורש של עץ הסיומות של  $s$ .