

תרגיל 8 – String matching

תאריך הגשה: 24/1/16

הוראות הגשה:

יש להגיש את התרגיל דרך האתר – submit.cs.biu.ac.il

יש לציין בתחלת הקובץ בהערה שם ות.ז.

יש להגיש קובץ ייחד בשם targil8_c.cpp (למי שmagish ב-C++) או targil8_c_.cpp (למי שmagish ב-C++).

תיאור המשימה:

עליכם לכתוב תוכנית המקבלת בשורת הפקודה מספר n ואחריו מחרוזת S שנקרה לה. לדוגמא, אם לקובץ הריצה שלכם קוראים a.out ומריצים את הפקודה:

```
>> a.out 5 aabAgFumy
```

אז $n=5$ ו- $S=aabAgFumy$.

התוכנית תעבור על כל תת-מחרוזות של S באורך n ותדפיס את כל תת-המחרוזות המופיעות **בדיקות פעמיים**. אם לא קיימת תת-מחרוזת כזו, לא תדפיסו כלום.

כל הדפסה כזו תהיה בשורה נפרדת. בין המחרוזות והמספרים יהיה רווח יחיד.

דוגמא של קלט:

```
>> a.out 3 the_rain_in_spain_stays_mainly_on_the_drain
```

הפלט (כמובן ללא מה שכותב באפור):

```
n_s 16 10 // the_rain_in_spain_stays_mainly_on_the_drain
the 34 0 // the_rain_in_spain_stays_mainly_on_the_drain
he_ 35 1 // the_rain_in_spain_stays_mainly_on_the_drain
rai 39 4 // the_rain_in_spain_stays_mainly_on_the_drain
```

שימוש לב שתתי המחרוזות "ain" ו- "in" לא יופיעו בפלט מכיוון שהן מופיעות יותר מפעם אחת!

על מנת שסדר ההדפסה יהיה זהה עבור כלם, עליכם **למיין** (בSİוכיות נורמלית...) את הפלט שלהם לפי המספר הראשוני (בסדר עולה), כלומר לפי מיקום המופיע השני של תת-מחרוזת מסוימת.

לדוגמא - אם אתם עוברים על המחרוזת ובודקים כל תת-מחרוזת, ואתם מקבלים את התוצאה הבאה:

```
the 34 0
he_ 35 1
rai 39 4
n_s 16 10
```

از עליים להדפיס את תת המחרוזות בצורה ממויינית בסדר הבא:

```
n_s 16 10
the 34 0
he_ 35 1
rai 39 4
```

מימוש באמצעות ריבין קראפ

- עליים להיעזר בפונקציית Hash של מחרוזת, שנסמך HS המקיים שnitן לחשב את $(a_n \dots a_2 a_1 a_0)_{HS}$ בזמן $O(1)$.
- עליים למשתמש בטבלת Hash, אליה תמפו ערכי Hash של תת מחרוזות באורך ℓ של S . אTEM רשאים למשתמש את טבלת ה-Hash בכל דרך שלמדתם, אך מומלץ למשתמש Cuckoo Hashing יותר תאים ממספר האיברים שמכניסים שהוא קל ונוח למימוש. כדי להקצות כ-50% יותר תאים ממספר האיברים שמכניסים לטבלת ה-Hash.
- המלצה נוספת – פונקציית Hash מוצלחת היא $(\sum a_i p^i) \% P = (a_1 a_2 \dots a_n)_{HS}$ כאשר d מספר ראשוני קטן ו- P מספר ראשוני גדול (לדוגמא, 999997 ב- $O(1)$).
- * כיצד תחשבו את hs_i מ- $hs_{i+1} \dots hs_0$ ב- $O(1)$? שימו לב שהפונקציה sow (העלאה בחזקה) לא עובדת ב- $O(1)$!

דרישות סיבוכיות:

זמן (ללא הרהיטס): $O(len(S))$.

מקום: $O(len(S))$.

הערה: מותר להשתמש בספריות `string.h`, `math.h` ובחלקת `vector` ב-`STL` (למי שmagic ב-`C++`).

מזכורת – אלגוריתם ריבין-קראפ:

בහינתן מחרוזת לחיפוש T ובמבנה P_0, \dots, P_k בגודל ℓ תוים, האלגוריתם מוצא את כל ההפניות של התבניות ב- T בזמן $(k+1)len(T)$. האלגוריתם משתמש בפונקציית hash מיוחדת HS המקיים שnitן לחשב את $(a_n \dots a_1 a_0)_{HS}$ בזמן $O(1)$ פעולות.

האלגוריתם בקצרה:

אם $(T[i:i+\ell-1] = P_k)$ פיתרון וסיימנו. אחרת:

בונים טבלת hash H בגודל k ולכל $i \leq k-1$ מבצעים $H[HS(P_i)] = hs_i$.
נסמן $(T[i:i+1] \dots T[i+\ell-1]) = hs_0$.

נחשב את hs_0 . אם $[hs_0]H$ מוגדר ומיצאנו התאמה ולאחר בדיקה ש- $P_{H[hs_0]}$ אכן מופיע ב- T (במקום המצוופה), נדפס את $[hs_0]H$.

כעת, עבור $(T[i:i+\ell-1] < 0)$ נבצע: חשב את hs_i מ- $hs_{i-1} \dots hs_0$. אם $[hs_i]H$ מוגדר מיצאנו התאמה ולאחר הבדיקה ש- $P_{H[hs_i]}$ אכן מופיע ב- T , נדפס את $[hs_i]H$.

דוגמאות קלט פלט:

```
>> a.out 4 AAAAaaaaaa
AAAA 1 0
aaaa 6 5
>> a.out 5 AAAAaaaaaa
>> a.out 4 SheSellsSeaShellsOnTheSeaShore
ells 13 4
heSe 20 1
SeaS 22 8
eaSh 23 9
```

```
>>> a.out 8 123456789101112
>>> a.out 2 123456789101112
11 12 11
12 13 0
>>> a.out 3 aaaAAAaaaAAaaa
aaA 7 1
aAA 8 2
AAA 9 3
AAa 10 4
Aaa 11 5
```