



# Keresés sztringben csúszóablakkal

Ciuciu-Kiss Jenifer Tabita  
Programtervező informatikus BSc

2018. január 22.

Az EFOP-3.6.3-VEKOP-16-2017-00002. számú projektben elvégzett szakmai feladat az Európai Unió támogatásával, az Európai Szociális Alap társfinanszírozásával valósult meg.

## Tartalomjegyzék

<b>1. Bevezetés</b>	<b>3</b>
<b>2. Keresés egy memóri felhasználásával</b>	<b>3</b>
2.1. Keresési fázis . . . . .	3
2.2. Tétel . . . . .	4
2.2.1. Bizonyítás . . . . .	5
2.3. A keresési fázis futási ideje . . . . .	6
2.4. Tétel . . . . .	6
2.4.1. Bizonyítás . . . . .	6
2.5. Következmény . . . . .	8
2.5.1. Bizonyítás . . . . .	9
<b>3. Keresés több memória felhasználásával</b>	<b>9</b>
3.1. Keresési fázis . . . . .	10
3.2. Lemma . . . . .	10
3.3. Lemma . . . . .	10
3.4. Lemma . . . . .	10
3.5. Lemma . . . . .	11
3.6. Az algoritmus . . . . .	12
3.7. Tétel . . . . .	12
3.7.1. Bizonyítás . . . . .	12
3.8. A keresési fázis komplexitása . . . . .	12
3.9. Következmény . . . . .	13
3.9.1. Bizonyítás . . . . .	13
3.10. Lemma . . . . .	15
3.10.1. Bizonyítás . . . . .	15
3.11. Lemma . . . . .	16
3.11.1. Bizonyítás . . . . .	16
3.12. Tétel . . . . .	16
3.12.1. Bizonyítás . . . . .	17
3.13. Következmény . . . . .	17
3.13.1. Bizonyítás . . . . .	17
<b>4. Keresés szótárban</b>	<b>17</b>
4.1. Az algoritmus . . . . .	19
4.2. Tétel . . . . .	20
<b>5. Összefoglalás</b>	<b>21</b>

**6. Források**

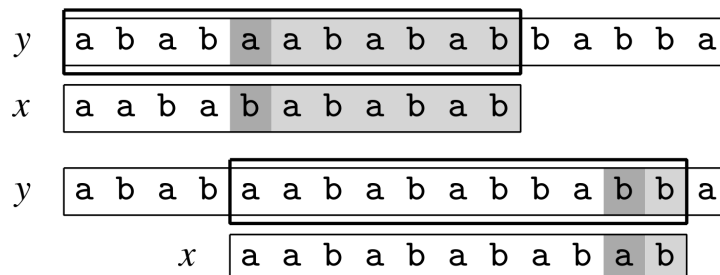
**22**

## 1. Bevezetés

Azt vizsgáljuk, hogy hogyan kell karakterláncokat keresni egy hosszabb szövegben. Erre több megoldást is láthatunk majd, melyek egyre több memóriafelhasználással ugyan, de egyre gyorsabban működnek.

## 2. Keresés egy memóri felhasználásával

Nézzük meg, hogy kell egy szövegben karakterláncot keresni hatékonyabban a naív, mindent mindennel összehasonlító algoritmusnál. A keresés során megjegyezzük legalább egy információt az előző egyezésről, jelen esetben az utolsó olyan szuffixet a sztringben ami a szövegben is előfordult. Ezt a technikát "faktor memorizációnak" nevezzük, ami kiterjeszti a prefix memorizációs implementációt, amit a szakirodalom WL-Memoryless-Suffix-Search -nek nevez. Ennek a keresésnek szüksége van egy konstans méretű extra memóriára a feldolgozandó szövegtől függően, hasonlóan az egyszerű Memoryless-Suffix-Search. A keresés működése azonban így már nem kvadrátikus, hanem  $2n$ -es, ahol  $n$  jelöli annak a szövegnek a hosszát amelyben keresünk.



1. ábra. Látható, hogy a 10-es pozícióban végzett kísérlet során felismertük az *ababab* szuffixet. Ezután eltoljuk az ablakot 4 pozícióval, és megjegyezzük, hogy az *ababababab* szuffix a szövegben 4 periódusú. Így az  $y$  tényező *aababab* értéke megegyezik az  $x$ -hez igazított tényező értékével. A következő lépésben felismerjük a  $b$  utótagot, az  $y[9] = a$  és az  $y[13] = b$  azt mutatják, hogy a szövegnek nincs 4-es periódusa. Tehát az *ababababab* szuffixe a 4-et ismeri fel periódusként, de az  $y[9] = a$  és az  $y[13] = b$  azt mutatják, hogy nem lehet őket egyidejűleg igazítani. Ez ahhoz vezet, hogy az eltolás  $|ababab| - |b| = 5$  hosszú lesz.

### 2.1. Keresési fázis

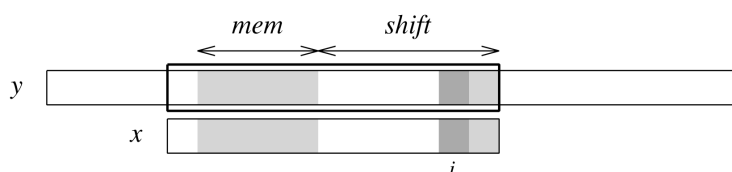
A Memoryless-Suffix-Search mindegyik eltolás után eldobja az korábbi kísérletek során összegyűjtött információkat. A jelenlegi algoritmus működését úgy gyorsítjuk, hogy megjegyezzük az  $x$  egy szuffixének utolsó előfordulását az  $y$  szövegben. Ennek az információnak a memorizációja több előnnyel is jár:

- lehetővé teszi a már biztosan nem jó tényezők átugrását
- lehetővé teszi az eltolások meghosszabítását

Az ebben az esszében ismertetett algoritmus nagy részben kihasználja ezeket a lehetőségeket. Az egyetlen tényező memorizálása egy pontosan adott szituációban lehetőséget. Az eltolás hosszának növelését a szakirodalom turbo-shiftnek nevezi, magyarosítva turbó-eltolásnak nevezik.

A következőkben részletesebben ismertetve lesz maga az algoritmus. Az általános szituáció egy  $T$  kísérlet során a keresési fázisban Turbó-suffix-keresés során bemutatásra került az 1. Az előző kísérlet legyen  $T'$ , a  $j'$  pozíción, a  $z'$  suffix lett felismerve a szövegben, és az eltolás hossza a következő lesz:  $d' = \text{bestfact}(m - 1 - |z'|, y[j' - |z'|])$ .

A  $T$  próbálkozás során a  $j = j' + d'$  pozícióban egy ugrás  $z'$  elem felett az  $y$  sztringben lehetséges, ha a  $d'$  hosszú suffix felismerhető az  $y$  sztringben. Ebben a esetben, felesleges összehasonlítani az  $z'$  elemet a a szövegből a keresett karakterláncsal, mert az eltolás definíciója alapján azok biztosan meg fognak egyezni. A turbó-eltolás használható, ha a  $z$  felismert suffix az aktuális próbálkozásban rövidebb, mint  $z'$ . Tehát a turbó-eltolás hossza:  $|z'| - |z|$ .



2. ábra. Az következő változók:  $i$ ,  $mem$  és  $shift$ , az instrukciók végrehajtása során a Turbó-Suffix-Keresés nevű, algoritmus 10-es sorában lévő adatok alapján láthatóak. A világos szürke rész az egyezést jelöli, a sötét szürke rész pedig azt amelyik nem egyezik.

Abban az esetben amikor a turbó-eltolás nagyobb mint a  $\text{bestfact}(m - |z|, y[j - |z|])$ , megjegyezzük, hogy az eltolás a jelenlegi lépésben lehet nagyobb, mint  $|z|$ . Az adott tényező megjegyzése csak akkor lehet végleges, ha a  $\text{bestfact}$  általi eltolás már ki lett számolva. Az módszer helyessége az argumentumok periodicitásán alapul.

Most következik maga a Turbó-Suffix-Keresés algoritmus. Az algoritmusban hivatkozunk a korábban már definiált  $\text{bestfact}$  értékére, de a  $\text{bestfact}$  helyett a  $\text{goodsuff}$  is használható az eltolás hosszának kiszámítására.

## 2.2. Tétel

Az  $x$  és  $y$  sztringekre alkalmazott Turbó-Suffix-Keresés megtalálja  $x$  összes előfordulását  $y$ -ban

```

TURBO-SUFFIX-SEARCH( $x, m, y, n$ )
1   $shift \leftarrow 0$ 
2   $mem \leftarrow 0$ 
3   $j \leftarrow m - 1$ 
4  while  $j < n$  do
5       $i \leftarrow m - 1$ 
6      while  $i \geq 0$  and  $x[i] = y[j - m + 1 + i]$  do
7          if  $i = m - shift$  then
8               $i \leftarrow i - mem - 1$     ▷ Jump
9          else  $i \leftarrow i - 1$ 
10     OUTPUT-IF( $i < 0$ )
11     if  $i < 0$  then
12          $shift \leftarrow per(x)$ 
13          $mem \leftarrow m - shift$ 
14     else  $turbo \leftarrow mem - m + 1 + i$ 
15         if  $turbo \leq best\text{-}fact(i, y[j - m + 1 + i])$  then
16              $shift \leftarrow best\text{-}fact(i, y[j - m + 1 + i])$ 
17              $mem \leftarrow \min\{m - shift, m - i\}$ 
18         else  $shift \leftarrow \max\{turbo, m - 1 - i\}$ 
19              $mem \leftarrow 0$ 
20      $j \leftarrow j + shift$     ▷ Shift

```

### 2.2.1. Bizonyítás

A különbség a Memoryless-Suffix-Search és a Turbó-Suffix-Keresés között azon alapul, hogy az eltolás hosszát növeljük, és a már biztosan egyező részeket nem hasonlítjuk össze. Ez elég, hogy belássuk, hogy a 18-as sorban kiszámolt eltolás helyes. Először lássuk be azt, hogy a turbó-eltolás (turbo változóként jelölve) hossza helyes. Ekkor belátjuk, hogy az  $m - 1 - i$  hosszú eltolás szintén helyes. Jegyezzük meg, hogy a 18. sor akkor hajtódik végre, ha  $turbo > best\text{-}fact(i, y[j - m + 1 + i])$ , ami implikálja azt, hogy  $turbo > 1$ .

A  $mem$  változó értéke a felismert  $z'$  szuffix hosszát jelöli az előző,  $T'$  lépésben. A felismert  $z = x[i + 1 \dots m - 1]$  szuffix hossza a  $T$  lépésben  $m - 1 - i$ . A  $turbo$  változó hossza  $|z'| - |z|$ . Legyen  $a = x[i]$  az a karakter, amely megelőzi a  $z$  szuffixet a szövegben, és legyen  $b = y[j - m + 1 + i]$  az a karakter, amelyik megelőzi a megfelelő előfordulását  $z$ -nek a szövegben. Legyen  $u = x[m - d \dots i]$ , (tehát van egy  $z'u z \preceq_{szuff} x$ ). Amikor  $z$

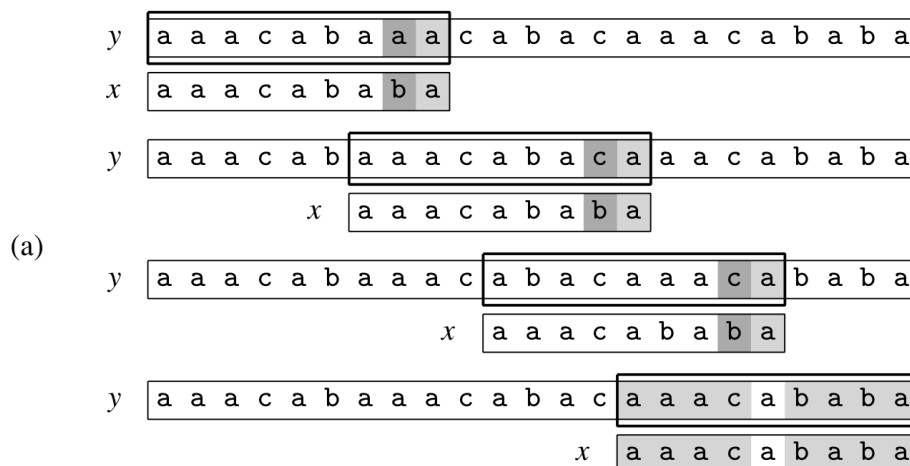
rövidebb mint  $z'$ , akkor  $az$  egy szuffix a  $z'$ -ben. Ebből következi, hogy az  $a$  és  $b$  karakterek előfordulnak a  $d = |uz|$  távolságra egymástól a szövegben. De a  $z'uz$  szuffix a sztringben is rendelkezik a  $d = |uz|$  periódussal (mert  $z'$  egy keret  $z'uz$ -ben), tehát az eltolások, amelyek rövidebbek, mint  $|z'| - |z| = turbo$  feleslegesen, hiszen nem vezetnek karakterkülönbözéshez. Tehát a *turbo* hosszú eltolás szabályos. Lásd 1.

Most lássuk be, hogy a  $|z| = m - 1 - i$  hosszú eltolás szabályos. Legyen  $l = bestfact(i b)$ . A *bestfact* definíciója alapján,  $x[i - l] \neq x[i]$ . Mivel a  $l$  egy periódusa  $z$ -nek, a  $x[i - l]$  és a  $x[i]$  karakterek nem lehetnek egyszerre a  $z$  előfordulásához igazítva az  $y$  tömbben. Ebből arra következtethetünk, hogy a  $|z| = m - 1 - i$  eltolás szabályos.

Két példa lesz majd látható a Turbó-Suffix-Keresés algoritmusára, az egyik a *bestfact*-ot, a másik a *goodsuff*-et használja a ??.

### 2.3. A keresési fázis futási ideje

Lássuk be, hogy a Turbó-Suffix-Keresés lineáris futási idejű a legrosszabb esetben.



### 2.4. Tétel

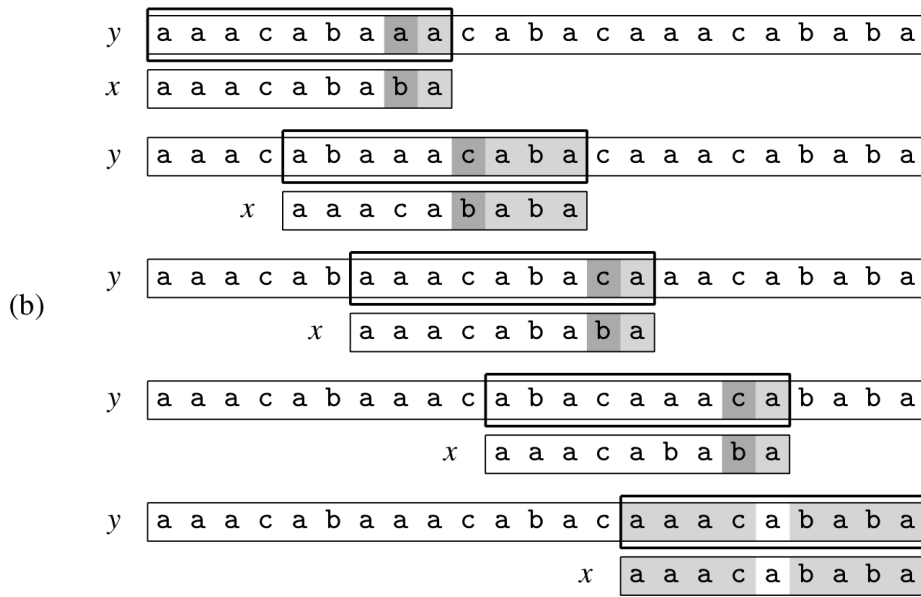
Az  $m$  hosszú  $x$  karakterlánc összes előfordulásának keresése egy  $n$  hosszú  $y$  szövegben a Turbó-Suffix-Keresés algoritmusával legrosszabb esetben  $2n$  karakterösszehasonlítást jelent.

#### 2.4.1. Bizonyítás

Felhasználva az előző bizonyítást, azt mondjuk, hogy a  $d$  hosszú eltolás a  $T$ -edik kísélet után rövid, ha  $2d < |z| + 1$ , és hosszú egyébként.

Különböztessünk meg három típusú próbálkozást:





3. ábra. Péda a Turbó-Suffix-Keresés két alkalmazására (a) *bestfact*-os használva, (b) *googsuff*-et használva a 18-as sorban az összehasonlításnál.

1. Olyan próbálkozás, amelyet egy ugrás követ,
2. Olyan próbálkozás, ami nem 1-es típusú és egy hosszú eltolás követi,
3. Olyan próbálkozás, ami nem 1-es típusú és egy rövid eltolás követi.

Az bizonyítás alapötlete az, hogy csökkentsük az összehasonlítások számát az eltolásokkal. Legyen a próbálkozások  $cost(T)$  a következő:

$$f := \begin{cases} 1, & \text{ha a } T \text{ próbálkozás 1-es típusú,} \\ 1 + |z|, & \text{ha a } T \text{ próbálkozás 2-es vagy 3-as típusú.} \end{cases}$$

Abban az esetben ha 1-es típusú egyezésünk van, a  $cost(T)$  függvény megadja, hogy mennyibe kerül a felesleges tesztelés. Az egyéb összehasonlítások költsége el lesz halasztva a következő összehasonlításra. Következésképpen az összes összehasonlítások száma a Turbó-Suffix-Keresés során egyenlő az összehasonlítások értékeinek összegével. Tehát bebizonyítottuk, hogy  $\sum_T cost(T) \leq 2 \sum_T d \leq 2n$ .

Egy 1-es típusú  $T_0$  próbálkozásra:  $cost(T_0) = 1 < 2d_0$ , ha  $d_0 \geq 1$ .

Egy 2-es típusú  $T_0$  próbálkozásra:  $cost(T_0) = |z_0| + 1 \leq 2d_0$ .

Már csak a 3-as típusú próbálkozás maradt hátra. Legyen egy 3-as típusú  $T_0$  próbálkozás a  $0$  pozíción az  $y$  szövegben. Ebben az esetben  $d_0 < |z_0|$ , és tudjuk, hogy  $d_0 = bestfact(m - |z_0|, y[j_0 - |z_0|])$ . Ez azt jelenti, hogy a következő  $T_1$  próbálkozás során a  $j_1$  és pozíción az  $y$  sztringben lehetséges egy turbó-eltolás

Nézzük a következő eseteket:

- a.  $|z_0| + d_0 \leq m$ . Ebből következik a turbó-eltolás definíciója alapján, hogy  $d_1 \geq |z_0| - |z_1|$ . Tehát:  $cost(T_0) = |z_0| + 1 \leq |z_1| + d_1 + 1 \leq d_0 + d_1$ .
- b.  $|z_0| + d_0 > m$ . Ebből következik a turbó-eltolás definíciója alapján, hogy  $|z_1| + d_0 + d_1 \geq m$ . Tehát:  $cost(T) \leq m \leq 2d_0 - 1 + d_1$

Mindig feltételezhetjük azt, hogy a b. eset következik be a  $T_1$ -es próbálkozásra, hiszen egy hosszú határt kapunk értékül a  $cost(T_0)$ -ra.

Amikor a  $T_1$  1-es típusú próbálkozásra tudjuk hogy  $cost(T_1) = 1$  és  $cost(T_0) + cost(T_1) \leq 2d_0 + d_1$ . Ez pedig jobb a vártnál.

Amikor a  $T_1$  2-es típusú próbálkozás vagy amikor  $|z_1| \leq d_1$  akkor  $cost(T_0) + cost(T_1) \leq 2d_0 + d_1$ .

Már csak az az eset maradt hátra, amikor  $T_1$  3-as típusú próbálkozás és  $|z_1| > d_1$ . Ez azt jelenti, hogy a  $T_0$  próbálkozás után  $d_1 = bestfact(m - |z_1|, y[j_1 - |z_1|])$ . A  $T_1$  próbálkozásnál használt argumentumot használjuk a következő  $T_2$  próbálkozásnál is. Az a. eset csak a  $T_2$  próbálkozásnál következhet be. Ez a következő eredményt adja  $cost(T_1) \leq d_1 + d_2$ . Végül  $cost(T_0) + cost(T_1) \leq 2d_0 + 2d_1 + d_2$ .

Az utolsó argumentum adja meg az indukciós bizonyítás lépését. Ha  $T_0, T_1, \dots, T_k$  mind 3-as típusú és  $|z_j| > d_j$  minden  $j = 0, 1, \dots, k$ , akkor

$$cost(T_0) + cost(T_1) + \dots + cost(T_k) \leq 2d_0 + 2d_1 + \dots + 2d_k + d_{k+1}.$$

Legyen  $T_{k'}$  az első próbálkozás a  $T_0$  próbálkozás után, tehát  $|z_{k'}| \leq d_{k'}$ . Ez a próbálkozás létezik, hiszen az ellentétes eset azt jelentené, hogy létezik egy végtelen ciklus az eseteknél az egyre rövidebb és rövidebb eltolások felé, ami lehetetlen. Tehát  $cost(T_0) + cost(T_1) + \dots + cost(T_{k'}) \leq 2d_0 + 2d_1 + \dots + 2d_{k'}$  és  $\sum_T cost(T) < 2 \sum_T d_T \leq 2n$ .

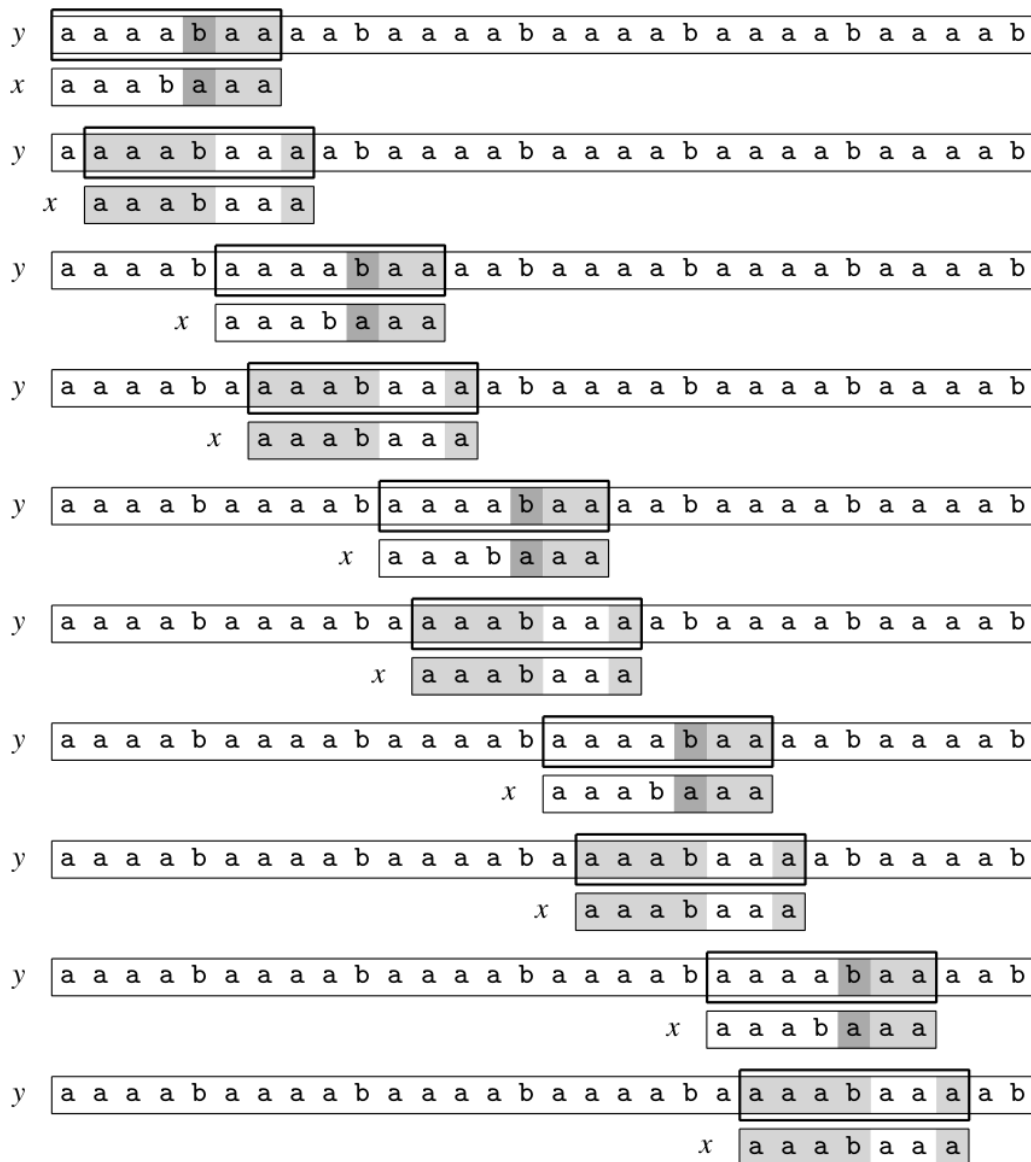
A határ a  $2n$  összehasonlítás a tétel esetén majdnem optimális. Nézzük a következő példát. Legyen  $x = a^k b a^k$  és  $y = (a^{k+1})^l$ , ( $k \geq 1$ ). Tudjuk hogy  $m = 2k + 1$  és  $n = l(k + 2)$ . Az  $a^{k+1}b$  ( $k + 2$  hosszú) első és az utolsó előfordulásától eltekintve az  $y$  szövegben a Turbó-Suffix-Keresés  $2k + 2$  összehasonlítást hajt végre. Elsőre  $k + 2$ -t, utolsóra  $k$ -t. Tehát ezt kaptuk.

$$(l - 1)(2k + 2) = 2n \left( \frac{m+1}{m+3} \right) - m - 1$$

karakter összehasonlítás történik az algoritmus működése közben. a 4-án látható egy példa  $k = 3$  és  $l = 6$  esetre.

## 2.5. Következmény

A Turbó-Suffix-Keresés algoritmus megtalálja az összes előfordulását egy adott karakterláncnak egy  $n$  hosszú szövegben  $O(n)$  idő alatt, konstans extra tárhely felhasználásával.



4. ábra. Példa a legrosszabb esetre a Turbó-Suffix-Keresés használatával.

### 2.5.1. Bizonyítás

Egyértelműen következik az előző tételből.

## 3. Keresés több memória felhasználásával

Most egy még gyorsabb algoritmussal ismerkedhetünk meg. Ezt a gyorsulást úgy érhetjük el, hogy még több információt jegyünk meg a korábbi próbálkozások során. Ehhez  $O(m)$  konstans méretű extra memóriára van szükség, azonban a futási idő  $1.5n$ -re csökken.

Ennek az algoritmusnak a neve Memória-Suffix-Keresés. Az algoritmus megjegyzi a keresett karakterlánc összes suffixének összes előfordulását a szövegben, és ezt az információt együtt használja a *suffix táblával*, annak érdekében, hogy még inkább megnöveljük

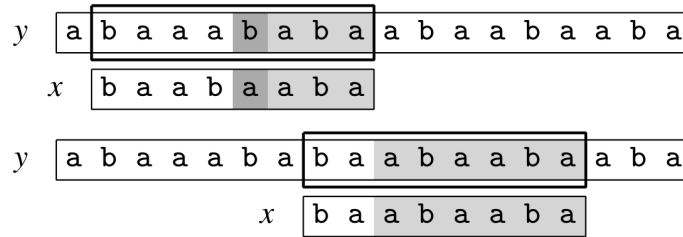
az eltolások hosszúságát. Ezek az eltolások a *bestfact* segítségével vannak kiszámolva, de kiszámíthatóak a *goodsuffix* felhasználásával is.

### 3.1. Keresési fázis

Itt kerül leírásra az algoritmus lényege. Minden próbálkozás után a  $j'$  pozíción az  $y$  szövegben, az  $x$  leghosszabb szuffixének hossza felismert a  $j'$  pozícióban,  $|j'|$  el van tárolva,  $|z'|$  pedig a táblában van eltárolva ( $S[j'] = |z'|$ ). A jelenlegi próbálkozás alatt a  $j$  pozíción az  $y$  szövegben, ha meg kell vizsgálnunk a  $|j'|$ ,  $j' < j$ , ahol a  $k = S[j']$ , tudjuk hogy  $y[j' - k + 1 \dots j'] \preceq_{szuff} x$ . Legyen  $i = m - 1 - j + j'$ . Elég tudni az  $x$  leghosszabb szuffixének a hosszát,  $s = suffix[i]$  az  $i$  pozíción végződve, hogy a legtöbb esetben következtetni tudjunk a próbálkozás eredményére.

### 3.2. Lemma

Amikor  $s \leq k$  és  $s = i + 1$ ,  $x$  egy előfordulása a jó  $j$  pozíción következik be az  $y$  szövegben. Tudjuk, hogy  $S[j] = m$  és hogy a  $per(x)$  hosszú eltolás szabályos.



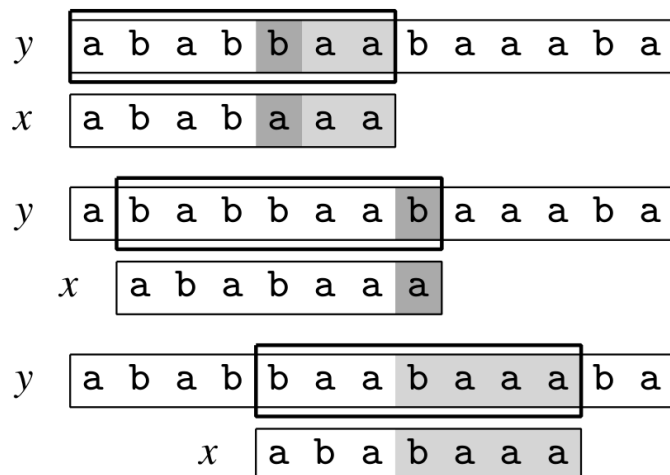
5. ábra. Egy próbálkozás alatt a 14-es pozíción észrevesszük, hogy a karakterenkénti összehasonlításakor az *abaaba* suffix 6 hosszú a sztringben. Elérünk a 8-as pozícióig az  $y$ -ban, ahol tudjuk, hogy a leghosszabb szuffix az  $x$ -en ezen a pozíción végződik és 3 hosszú. Tudjuk, hogy az  $x$  leghosszabb szuffixe az 1-es pozíción ér véget. Az első lemma alapján tudjuk, hogy ez az  $x$  egy előfordulását jelenti a 14-es pozíción.

### 3.3. Lemma

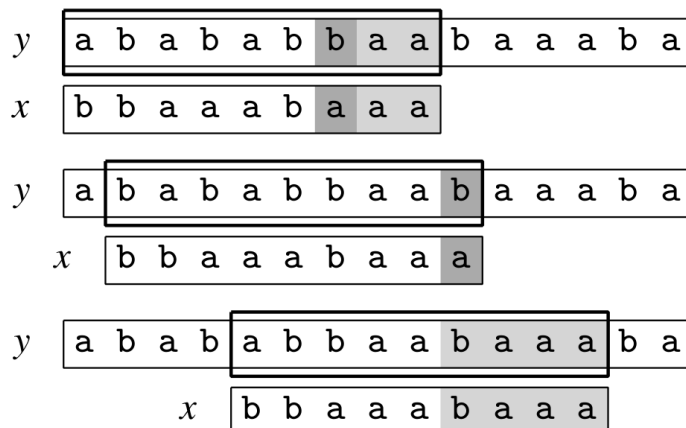
Amikor  $s \leq i$  és  $s < k$  akkor tudjuk hogy  $S[j] = m - 1 - i + s$  és legyen  $j' = j - m + 1 + i$ , ekkor a  $bestfact(i - s, y[j' - s])$  hosszú eltolás szabályos.

### 3.4. Lemma

Amikor  $k < s$ , akkor tudjuk, hogy  $S[i] = m - 1 - i + k$ , és legyen  $j = j - m + 1 + 1$ , ekkor tudjuk hogy a  $bestfact(i - k, y[j' - k])$ .



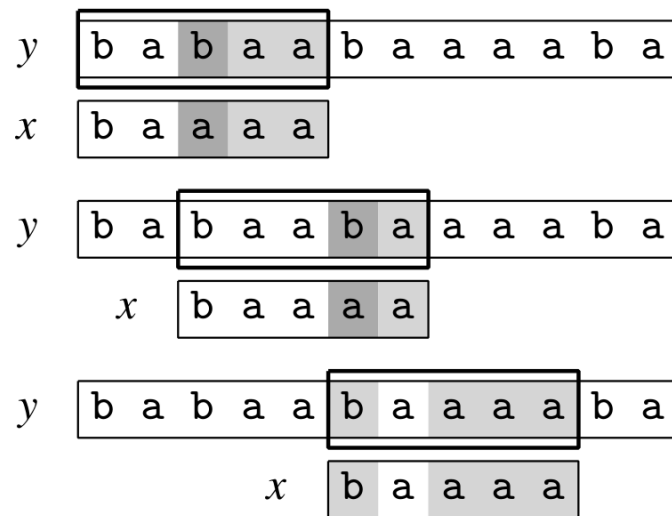
6. ábra. Amíg a 10-es pozíción próbálkozunk, észrevehettük, hogy a *baaa* a szöveg 4 hosszú szuffixe. Ekkor elmehetünk az  $y$  6-os pozíciójára, ahol tudjuk, hogy az  $x$  leghosszabb szuffixe itt végződik és 2 hosszú. Tehát az  $y[4 \dots 6]$  összehasonlítása felesleges, mert ott nem lesz egyezés.



7. ábra. Amíg a 12-es pozíción próbálkozunk, észrevehettük, hogy a sztring 4 hosszú szuffixe a 8-as pozíción van, ahol tudjuk, hogy a leghosszabb szuffix ami ezen a helyen végződik az 2 hosszú. Tehát egyetlen összehasonlítás nélkül is tudjuk, hogy  $x[2] = a$  és  $y[6] = b$ .

### 3.5. Lemma

Amikor  $k = s$ , akkor tudjuk, hogy  $x[i - s + 1 \dots m - 1] = y[j' - s + 1 \dots j]$ , és  $S[j] = m - 1 - i + s + |lcsuff(x[0 \dots i - s], y[j - m + 1 \dots j' - s])|$  ami azt jelenti, hogy  $j' = j - m + 1 + i$ .



8. ábra. Amíg a 9-es pozíció próbálkozunk, észrevehetjük, hogy a sztring 3 hosszú szuffixe a 6-os pozícióra került, ahol tudjuk, hogy a leghosszabb szuffix a sztringben ezen a pozíció végződik és 1 hosszú. Emellett tudjuk, hogy a leghosszabb szuffix a sztringben ami az 1-es pozíció végződik 1 hosszú. Tehát átgorthatjuk a  $y[6]$ -ot, mert összehasonlítás nélkül is tudjuk, hogy  $x[0] = y[5]$ .

### 3.6. Az algoritmus

Most pedig következzen maga az algoritmus 3.6. Jelen esetben a *best fact*-ot és a szuffix táblát használtuk. Az  $x$  szuffixeinek előfordulásainak memorizálásához a szövegben a szuffix táblában kerül tárolásra. Ennek a táblának az értékeire kezdetben mindenhol 0 prioritást állítunk be.

### 3.7. Tétel

A Memória-Szuffix-Keresés megtalálja egy  $x$  karakterlánc összes előfordulását egy  $y$  szövegben.

#### 3.7.1. Bizonyítás

Egyértelműen következik az első és az utolsó lemmából.

### 3.8. A keresési fázis komplexitása

Eddig folyamatosan azt vizsgáltuk, hogy a mekkora a tárigénye a Memória-Szuffix-Keresésnek futás közben.

```

MEMORY-SUFFIX-SEARCH( $x, m, y, n$ )
1  for  $j \leftarrow 0$  to  $n - 1$  do
2       $S[j] \leftarrow 0$ 
3   $j \leftarrow m - 1$ 
4  while  $j < n$  do
5       $i \leftarrow m - 1$ 
6      while  $i \geq 0$  do
7          if  $S[j - m + 1 + i] > 0$  then
8               $k \leftarrow S[j - m + 1 + i]$ 
9               $s \leftarrow \text{suff}[i]$ 
10             if  $s \neq k$  then
11                  $i \leftarrow i - \min\{s, k\}$ 
12                 break
13             else  $i \leftarrow i - k \triangleright$  Jump
14             elseif  $x[i] = y[j - m + 1 + i]$  then
15                  $i \leftarrow i - 1$ 
16             else break
17     OUTPUT-IF( $i < 0$ )
18     if  $i < 0$  then
19          $S[j] \leftarrow m$ 
20          $j \leftarrow j + \text{per}(x)$ 
21     else  $S[j] \leftarrow m - 1 - i$ 
22          $j \leftarrow j + \text{best-fact}(i, y[j - m + 1 + i])$ 

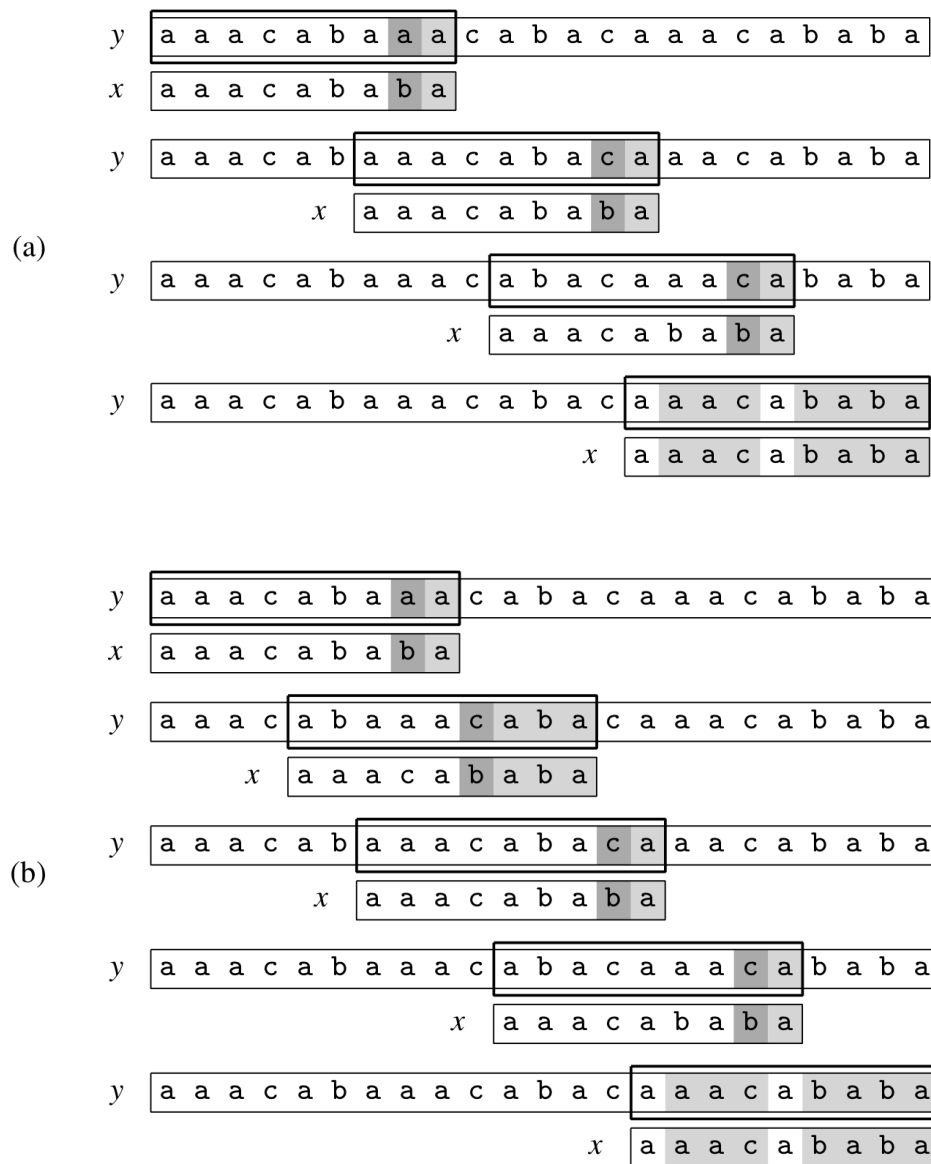
```

### 3.9. Következmény

Egy  $m$  hosszú  $x$  karakteránc keresése egy szövegben megoldható  $\mathbf{O}(m)$  extra memória felhasználásával.

#### 3.9.1. Bizonyítás

A tárhelyet a szuffix tábla memorizálásához használjuk, hogy implementáljuk a *best fact* funkciót vagy a *goodsuff* funkciót az  $S$  szuffix tábla létrehozásához. Az első 3 elem elhasznál  $\mathbf{O}(m)$  memóriát. Az  $S$  táblában nem jegyezzük meg csak a  $S[j - m + 1 \dots j]$ , mert csak ez a része a hasznos amikor a kereső ablak a jobboldali  $j$  pozíción van. Az  $S$



9. ábra. Két példa a Memória-Suffix-Keresés algoritmusára. (a) *bestfact* használatával. Ebben az esetben a 13 az összehasonlítások száma a keresett karakterlánc karakterei és a szöveg karakterei között. (b) a *goodsuff* használatával. Ebben az esetben ugyanez a szám 17.

tábla menedzseléséhez úgy tekintünk rá mint egy ciklikus listára, amiben csak a hasznos elemeket tároljuk, azaz  $S[l]$  elemet, tehát ez így még mindig  $O(m)$  tárigényű.

Most pedig lássuk be, hogy a futási ideje  $O(n)$ . Ez egész pontosan  $1.5n$  karakterösszehasonlítást jelent ahhoz, hogy megtaláljuk egy  $x$  karakterlánc összes előfordulását egy  $y$  szövegben.

Jegyezzük meg először is, hogy a keresési fázis alatt egy karakter előfordulása az  $y$  szövegben csak egyszer lesz összehasonlítva bármivel, egy karaktert a legtöbb esetben nem hasonlítunk kétszer. Ez tehát legrosszabb esetben  $n$  darab összehasonlítást jelent.



Csak azokat a karaktereket hasonlíthatjuk össze kétszer, amelyek egy előzőleg egy nem egyezést mutattak.

Nézzük a következő lemmát.

### 3.10. Lemma

Egy próbálkozás során a Memória-Suffix-Keresés algoritmusában ha  $k$  pozitív összehasonlítást már elvégeztünk, akkor a következő eltolás legább  $k$  hosszú lesz.



10. ábra. Egy próbálkozás alatt a  $j$  pozíción a szövegben észrevehettük, az 1 hosszú szuffixet a sztringben, ekkor pedig eltolhatjuk a 6-os pozícióba. Vegyük észre, hogy ekkor a szuffix 3 hosszú, tehát eltolhatjuk 4 pozícióval a keresőablakot. Ekkor láthatjuk, hogy itt egy 2 hosszú szuffix van, tehát eltolhatjuk az ablakot 5 pozícióval. A  $j + 15$  próbálkozás során észrevehettük, hogy a 19 hosszú szuffix egyezik, mindez annak árán, hogy a  $y[j + 8]$ ,  $y[j + 3]$ ,  $y[j - 1]$  pozíción lévő elemeket kétszer is összehasonlítottuk. Az eltolás ami ezt a próbálkozást követi nem lehet rövidebb vagy egyenlő mint 3 a negyedik lemma alapján. Valóban a szuffix  $aabaaabaaaabaaaaaa$  a szövegben nem lehet kisebb vagy egyenlő periódusú mint 3.

#### 3.10.1. Bizonyítás

Legyen  $T$  egy kísérlet a Memória-Suffix-Keresés algoritmus során. Legyen  $k$  az a karakter, amelyet más összehasonlítottunk egyszer és újra hasonlítottunk. Legyen a szöveg:

$$b_0 v_0 b_1 v_1 \dots b_k u_k v_k$$

a megmaradt tényező a  $T$  kísérlet során.

$\rightarrow b_0$  az a karakter amelyik nem egyezést okoz  $T$  során,  $\rightarrow b_0 v_0 b_1 v_1 \dots b_k u_k v_k$  egy szuffixe  $x$ -nek,  $\rightarrow a b_l$ , ( $1 \leq l \leq k$ ) karakterek azok a karakterek amelyek újra lettek próbálva  $T$  során és egyezést mutattak,  $\rightarrow a u_l$ , ( $1 \leq l \leq k$ ) azok a szuffixek (valószínűleg üresek), amelyek azok a karaktereket tartalmazzák amelyek össze lettek hasonlítva egyezően, és  $b_l$  az amelyek nem egyeztek, ezek a tényezők át lesznek ugorva a  $T$  kísérlet

során,  $\rightarrow a v_l$ , ( $1 \leq l \leq k$ ) tényezők azok amelyek egyezően lettek összehasonlítva elsőre  $T$  során.

Definíció alapján  $b_l u_l$ , ( $1 \leq l \leq k$ ) nem szuffixe  $x$ -nek.

A bizonyítás az összegésen alapszik. Feltételezzük, hogy  $d$  eltolás alkalmazva lett csak a  $T$  kísérlet után, és ez rövidebb mint  $k$ . Ekkor legyen  $w$  az  $x$ -nek az a szuffixe amelyik  $d$  hosszú. Ekkor a definícióból következik, hogy

$$v_0 b_1 u_1 v_1 b_2 u_2 v_2 \dots b_k u_k v_k w$$

egy szuffixe  $x$ -nek úgy hogy  $d = |w|$ .

Ekkor két különböző indexre  $l' \neq l''$ ,  $u_{l'}$  és  $u_{l''}$  ugyanazzal a  $w$  tényezővel vannak határolva, és ez maximum  $k - 1$  pozíción lehetséges. Ez azt jelenti, hogy  $b_{l'} u_{l'} = b_{l''} u_{l''}$ . Az eltolás, amely alkalmazva lett a két kísérletre  $b_{l'}$  és  $b_{l''}$ . Ekkor létezik egy olyan  $l$  index,  $l < k$  és  $b_l u_l = b_k u_k$ , ami ellentmond annak a ténynek, hogy az  $x$  karakterlánc már határolva volt a  $T$  kísérlet előtt.

Ez azt jelenti, hogy a  $T$  kísérlet során alkalmazott eltolás legalább  $k$  hosszú.

### 3.11. Lemma

A Memória-Suffix-Keresés algoritmus legfeljebb  $n/2$  olyan karaktert hasonlít újra, amelyet már egyszer összehasonított.

#### 3.11.1. Bizonyítás

Csoportosítsuk a próbálkozásokat. Két próbálkozás azonos csoportba kerül, ha van olyan karakter amelyet mindketten összehasonlítanak a szövegben. Egy  $p$  csomag amelyben  $k$  darab egyező újraösszehasonlítás volt és ebben a csomagban legyen legalább  $k + 1$  próbálkozás. Ezek alapján van  $k$  darab próbálkozásunk amelyek egy legalább 1 hosszú eltolással járnak, illetve van egy eltolásunk, ahol az eltolás legalább  $k$  az 5-ös lemma alapján. Tehát az együttes eltolás összesen legalább  $2k$ .

Ebből következik hogy az összes csomag együttes eltolása a Memória-Suffix-Keresés közben nem több mint  $n$ . Tehát az újraösszehasonlítások száma nem több mint  $n/2$ .

### 3.12. Tétel

Egy  $n$  hosszú  $y$  szövegben való keresés közben a Memória-Suffix-Keresés algoritmus legrosszabb esetben  $1.5n$  összehasonlítást tartalmaz a szöveg és a keresett karakterlánc betűi között.

### 3.12.1. Bizonyítás

Egyértelműen következik az előző lemmából és abból a tényből, hogy legfeljebb  $n$  egyező összehasonlítás lehet.

## 3.13. Következmény

A Memória-Suffix-Keresés algoritmus egy  $m$  hosszú karakterlánc összes előfordulását megtalálja egy  $m$  hosszú szövegben  $\mathbf{O}(n)$  idő alatt és  $\mathbf{O}(m)$  extra memória felhasználásával.

### 3.13.1. Bizonyítás

Egyértelműen következik az ezt meglőző és az azelőtti tételből.

## 4. Keresés szótárban

A csúszó ablak technikával megoldható, hogy keressük egy karakterlánc összes előfordulását egy  $k$  karakterláncból álló  $X = x_0, x_1, \dots, x_{k-1}$   $y$  szövegben. Ez most legyen  $m'$  és  $m''$  rendre a leghosszabb és a legrövidebb karakter hossza  $X$ -ben.

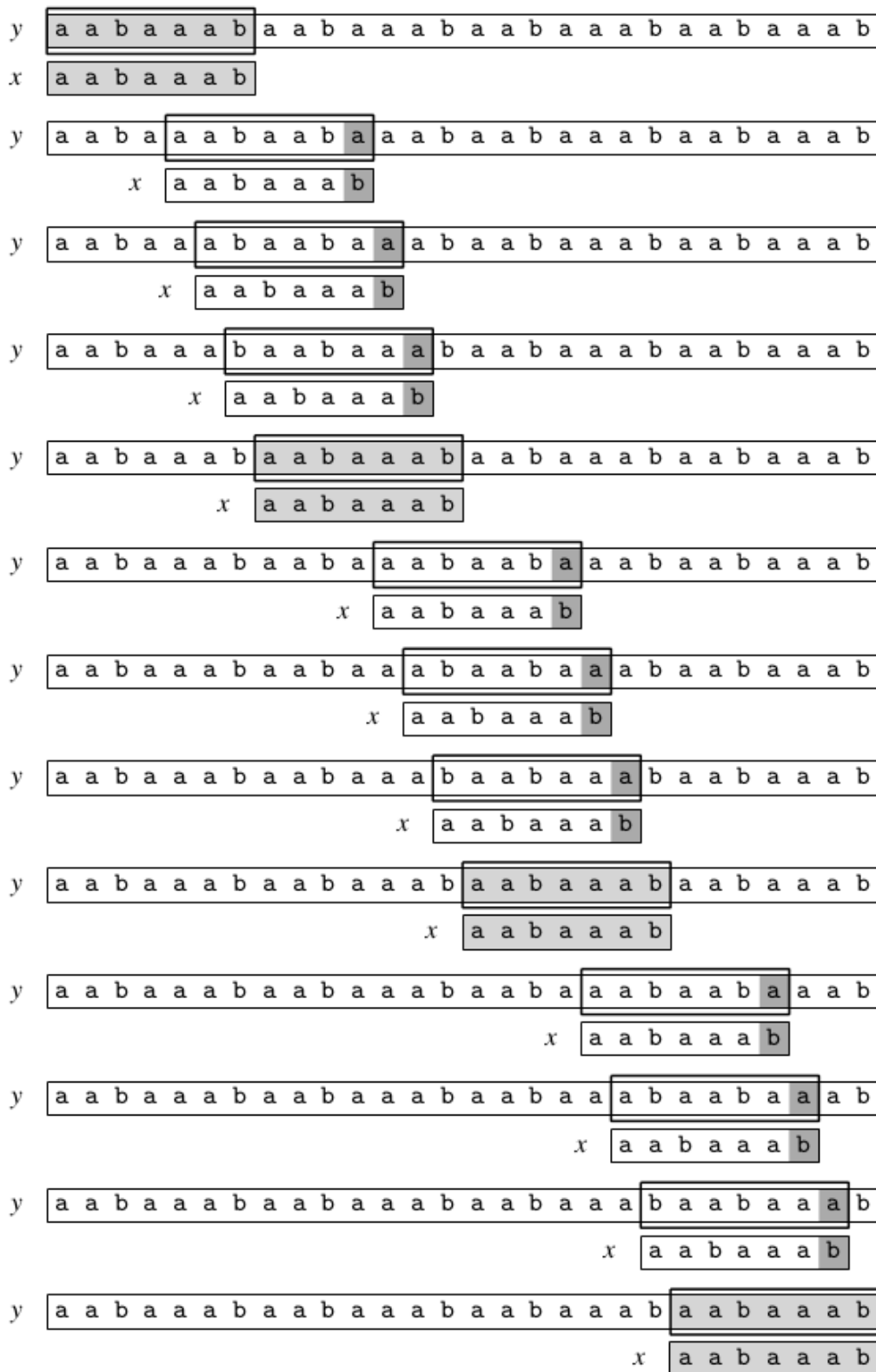
A szöveg beolvasásának hosszát az határozza meg, hogy mekkora a leghosszabb tényezője  $X$ -nek, ami egy suffixe lesz a csúszó ablakunknak. Ilyen módon hosszabbítsa meg a suffixet a beolvasott ablakhoz képest az előző módszer alapján. Ez lehetőséget ad arra, hogy több információt gyűjtsön a szövegről, és gyakran hosszabb eltolódáshoz vezet. Így az  $y$  szöveg beolvasása során az algoritmus elegendő információt fog tartalmazni az keresett karakterláncok előfordulási helyeinek észleléséhez.

Az algoritmus elsődleges célja az  $X$  karakterláncainak felderítése az  $m''$  hosszúságú csúszóablakban. Az alapelv az, hogy minden próbálkozás során meghatározzuk az  $x$  karakterláncainak prefixeit, amit a csúszóablak suffixei lesznek. Ezzel egyidejűleg figyeljük, hogy az  $X$  elemeinek előfordulásai az ablakban, és megjegyeztük azt az eltolódást aminek nem tolhatunk el hosszabban, ez legyen  $m'$ .

Most ismertessük azt a technikát amelyet ennek megvalósítására használunk. Legyen  $X^\sim$  az a karakterláncokból álló szöveg, amelyik az  $X$  egyes karakterláncainak megfordítását tartalmazza. Legyen  $N$  egy olyan determinisztikus automata, ami felismeri az  $X^\sim$  sztringek utótagjait. A hozzá rendelt átmeneti függvény legyen  $\delta$ . Az automata a következő nyelvet ismeri fel

$$\text{Suff}(X^\sim) = \{v \in A^* : uv = x^\sim, u \in A^*, x \in X\}.$$

Vagyis  $N$  determinisztikusan felismeri a prefixeit  $X$  elemeinek. Minden terminális állapotára az automatának legyen:



11. ábra. Példa a tételre,  $x = aabaaab$  és  $y = (aabaaab)^4$ . A keresett karakterlánk 7 hosszú, egy 28 hosszú szövegben. Az összehasonlítások száma 37. Az  $x$  utolsó 3 előfordulása esetén a Memória-Suffix-Keresés algoritmus 10 összehasonlítást végez.

$$output[q] = i : 0 \leq i \leq k - 1 \wedge \bar{\delta}(q_0, x_i^{\sim}).$$

ahol  $\delta$  az automata  $\delta$  átmeneti függvényeinek kiterjesztése, a  $q_0$  pedig a kezdőállapot.

Egy kísérlet a  $j$  pozíción az  $y$  szövegben az  $y$  karaktereinek elemzéséből áll a bal oldali  $y[j]$ -től az  $N$  automata által. Minden alkalommal, amikor a  $q$  állapot elér egy  $y[j']$  karakterig, ellenőrizzük hogy az  $output[q]$  üres e. Ha ez a helyzet, akkor az  $x_i$  sztring megjelenik az  $y$ -ban a  $j'$  pozíción, amikor

$$i \in output[q] \wedge j - j' + 1 = |x_i|.$$

Emellett ha a  $q$  állapot terminálóállapot, akkor az eltolás hossza nem lehet nagyobb, mint  $m' - (j - j' + 1)$ , feltéve, hogy ez a mennyiség pozitív. Így számíthatunk az érvényes eltolások minimális  $d$  hosszúságára. Egy kísérlet akkor ér véget, amikor az aktuális állapottól a következő karakteren keresztül nincs több átmenet.

#### 4.1. Az algoritmus

MULTIPLE-SUFFIX-SEARCH( $X, m', y, n$ )

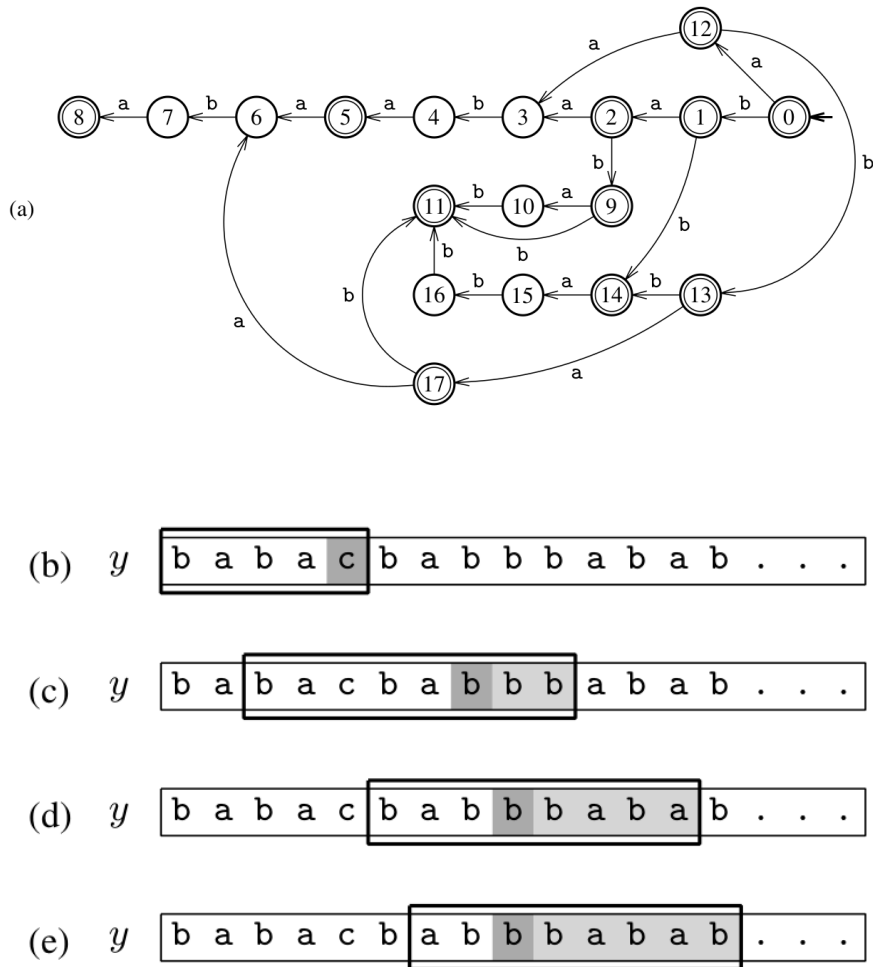
```

1  let  $N$  be a (deterministic) automaton
   accepting the suffixes of the reverse strings of  $X$ 
2   $j \leftarrow m' - 1$ 
3  while  $j \leq n$  do
4       $q \leftarrow initial[N]$ 
5       $j' \leftarrow j$ 
6       $d \leftarrow m'$ 
7      while  $j' \geq 0$  and TARGET( $q, y[j']$ )  $\neq$  NIL do
8           $q \leftarrow TARGET(q, y[j'])$ 
9          if terminal[ $q$ ] then
10             for each  $i \in output[q]$  do
11                 OUTPUT-IF( $|x_i| = j - j' + 1$ )
12             if  $m' - j + j' - 1 > 0$  then
13                  $d \leftarrow \min\{d, m' - j + j' - 1\}$ 
14             else  $d \leftarrow 1$ 
15              $j' \leftarrow j' - 1$ 
16      $j \leftarrow j + d$ 

```

### 4.2. Tétel

A Többszörös-Suffix-Keresés algoritmus felismeri az  $X$  karakterláncalmaz összes elemének összes előfordulását egy  $y$  szövegben.



12. ábra. A Többszörös-Suffix-Keresés algoritmus egy példán keresztül bemutatva.

## 5. Összefoglalás

Szó volt arról, hogy hogyan kell karakterláncokat keresnünk egy hosszabb szövegben. Erre több megoldást is láthattunk, melyek egyre több memóiafelhasználással ugyan, de egyre gyorsabban működtek.

Először csak egy segédmemóriát használtunk, ennek az algoritmusnak a neve Turbó-Suffix-Keresés. itt használtuk fel a legkevesebb extra tárhelyet, azonban ez az algoritmus a legrosszabb esetben  $2n$  karakterösszehasonlítást végez, ahol  $n$  annak a szövegnek a hossza, amelyben keresünk.

Ezt követően már több extra tárhelyet is igénybe vettünk, azonban így el tudtuk érni, hogy a legrosszabb esetben is csak  $1.5n$  összehasonlítást végezzünk. Ennek az algoritmusnak a neve Memória-Suffix-Keresés volt.

Ezt követően megnéztük, hogy hogyan kell egy több szót tartalmazó halmazban keresni egy adott karakterlánc előfordulását.

Mindegyik algoritmusra néztünk példát, elemeztük a működését.

## 6. Források

→ Maxime Crochemore, Christophe Hancart, Thierry Lecroq: Algorithms on strings, Cambridge University Press, Paris, 2001