

Közelítő illesztés és indexelés részstringekkel

Lehotay-Kéry Péter

lepuaai@inf.elte.hu

Ben Langmead diasora alapján

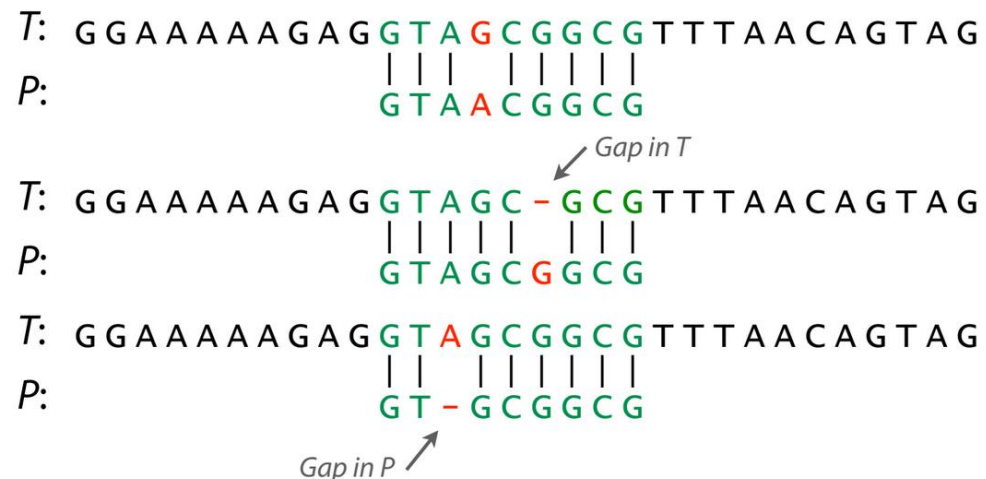
(www.langmead-lab.org/teaching-materials)

Szükség van közelítő illesztésre

- Szekvencia eltérések miatt
 - Szekvenálási hibák
 - Természetes variációk

Közelítő string illesztés

- Olyan helyeket keresünk, ahol P illeszkedik T-re, megengedve egy adott számú eltérést, vagy szerkesztést.
- Egy eltérés egy karakternyi helyettesítés.
- Egy szerkesztés lehet egy karakternyi helyettesítés, vagy hézag.



Hamming és szerkesztés távolság

- 2 azonos hosszú X és Y stringre, a hamming távolság a szükséges minimális számú egy-karakteres behelyettesítés, ahhoz hogy egyik stringet a másikká változtassuk.

X: G A G G T A G C G G C G T T T A A C
| | | | | | | | | | | | | | | |
Y: G T G G T A A C G G G G T T T A A C

Hamming distance = 3

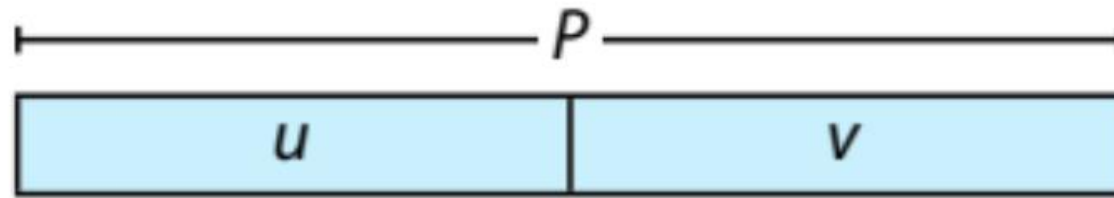
- Szerkesztés távolság a szükséges minimális számú szerkesztés, ahhoz, hogy egyik a másikká váljon.

X: T G G C C G C G C A A A A A C A G C
| | | | | | | | | | | | | | | |
Y: T G A C C G C G C A A A A - C A G C

Edit distance = 2

Közelítő string illesztés

- Hogyan tegyük Boyer-Moore és indexek által segített illesztéseket közelítővé?
- Osszuk P -t nem-üres, nem-átfedő részstringekre: u és v . Ha P előfordul T -ben 1 szerkesztéssel, u -nak, vagy v -nek pontosan egyeznie kell.

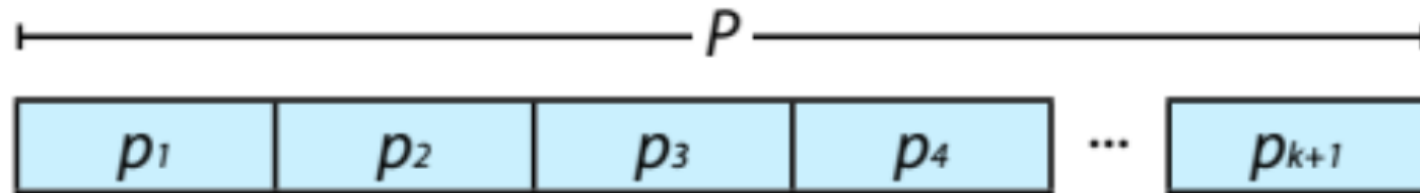


Either the edit goes here...

...or here. Can't go anywhere else!

Közelítő string illesztés

- Általánosabban: Legyen p_1, p_2, \dots, p_{k+1} P particionálása $k+1$ nem-átfedő, nem-üres stringekbe. Ha P előfordul T -ben maximum k szerkesztéssel, akkor p_1, p_2, \dots, p_{k+1} közül legalább egynek pontosan kell egyeznie.



- Használjuk a pontos illesztés algoritmust hogy pontos illeszkedést találjunk p_1, p_2, \dots, p_{k+1} -el. Keressünk hosszabb közelítő illeszkedést a pontos illeszkedés körül.

Előfeldolgozás

- Láttuk a naive algoritmust és Boyer-Moore-t. Ezek és más algoritmusok szétválaszthatók az alapján, milyen előfeldolgozást végeznek.
- Naiv nem végez előfeldolgozást.
- Boyer-Moore előre feldolgozza P mintát:
 - Ha P előre adott, kikereső táblákat építhetünk a rossz karakter és jó szuffix szabályokhoz.
 - Ha T1 később jön, használjuk az előre felépített táblákat, hogy összevessük P-t T1-el
 - T2, T3, ..-vel ugyan ez a helyzet.

Indexelés

- P-t, T-t vagy mindkettőt is feldolgozhatjuk előre. Amikor T-t dolgozzuk fel, azt mondjuk indexet készítünk, mint egy könyv indexe.
- Webes keresőmotorok is indexelést használnak:

Documents:

T[0] = "it is what it is"

T[1] = "what is it"

T[2] = "it is a banana"

http://en.wikipedia.org/wiki/Inverted_index

Index:

"a": [2]

"banana": [2]

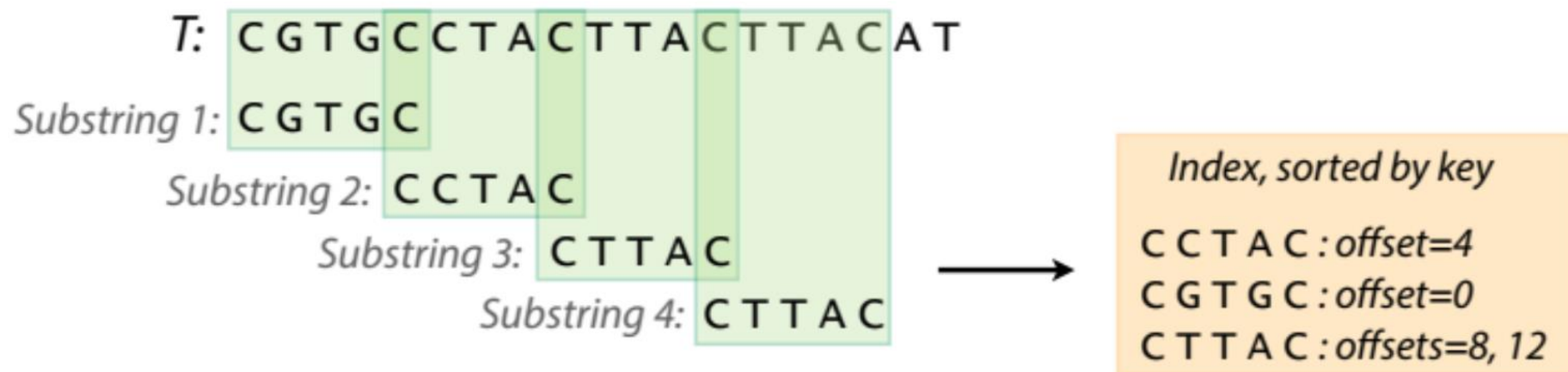
"is": [0, 1, 2]

"it": [0, 1, 2]

"what": [0, 1]

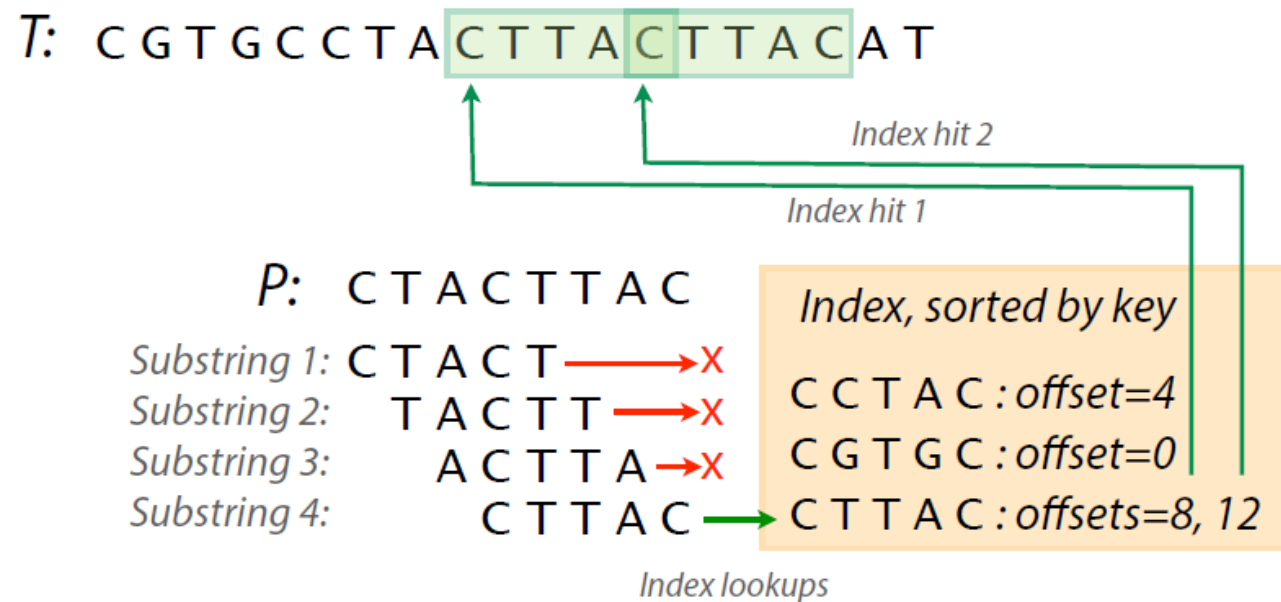
Részstring index

- T indexelése:
 - Kinyerjük a szekvenciákat T-ből (általában részstringeket), mutatókkal, hogy hol fordultak elő.
 - Rendezzük a darabokat és mutatókat egy adatszerkezetbe. Különböző szerkezetek lehetnek, mindegyik magába foglal valamilyen csoportosítást, rendezést.



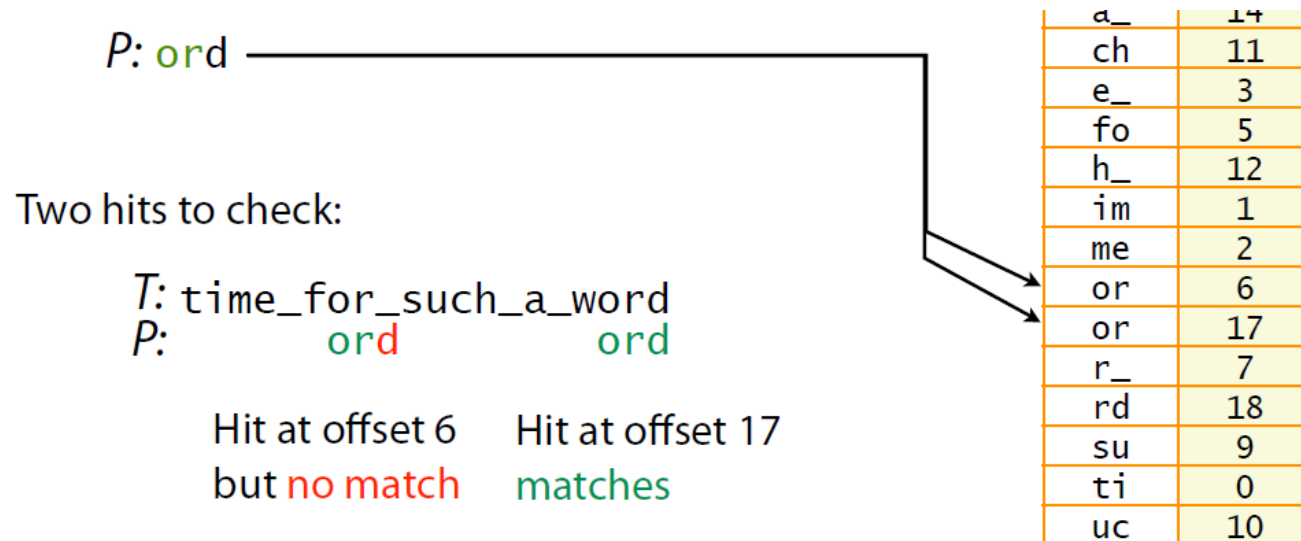
Részstring index

- Index lekérdezése P-vel:
 - Kinyerjük a részstringeket P-ből. Ezeket használjuk, hogy lekérdezzük az indexet. Index megmondja nekünk, hol fordulhatnak elő T-ben.
 - Minden előfordulásra, megnézzük, van-e illeszkedés ezen a környéken



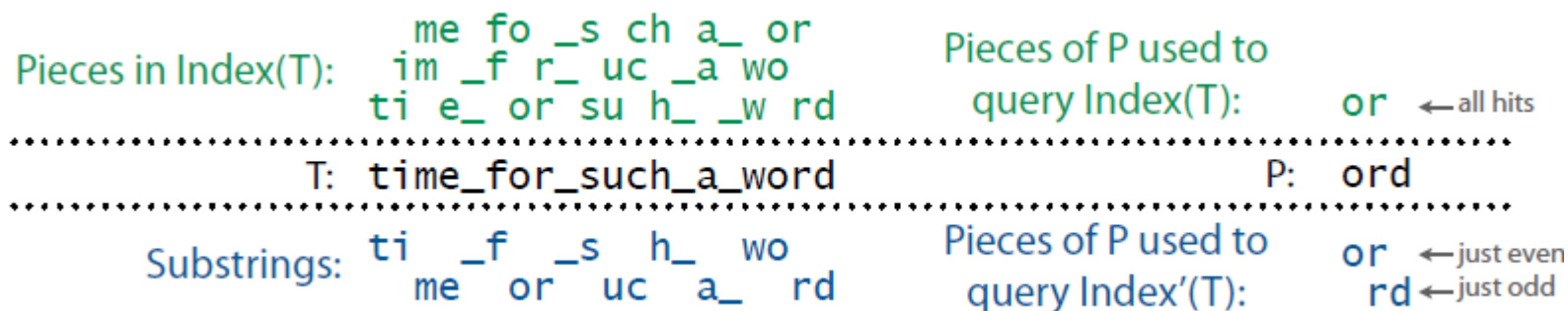
Részstring index

- Index T-hez: rendezett tábla, minden 2-hosszú részstringhez és előfordulásukhoz.
- Lekérdezéskor: kinyerjük a 2-hosszú prefixét P-nek, kikeressük őket az indexből, megvizsgáljuk a találatokat.



Részstring index: minden második részstring

- Ahelyett, hogy minden 2-hosszú részstringet kinyerjük, kihagyunk párat, mondjuk minden másodikat.
- T lekérdezésekor:
 - Kinyerjük 2 legszélső 2-hosszú részstringjét P-nek, megkeressük az indexben, kipróbálunk minden nevezőt.
 - Első kikeresés megtalál majd minden illeszkedést a páros előfordulásokra, második a páratlan előfordulásokra.



T: time_for_such_a_word

Substring	Offset
_f	4
_s	8
a_	14
h_	12
me	2
or	6
rd	18
ti	0
uc	10
wo	16

Részstring index: minden N-edik részstring

- Ahogy kivettük minden 2. részstringet, úgy kivehetjük minden 3., 4., stb
- Ha kivesszük minden N. részstringet, használnunk kell P első N részstringjét a lekérdezéshez.

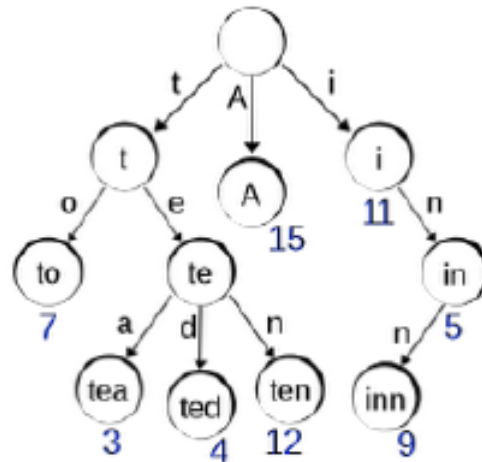
Indexelés: adatszerkezetek

- Indexek párosítják a részstringeket, vagy részszekvenciákat a találatokhoz.
- Több alkalmas adatszerkezet is létezik:

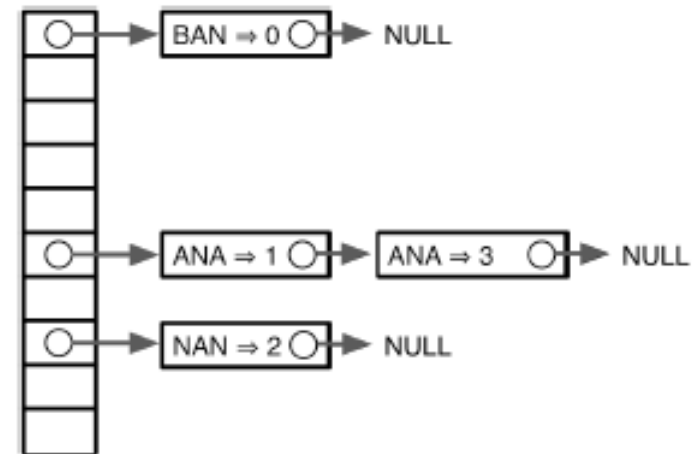
Sorted list:

_f	4
_s	8
a_	14
h_	12
me	2
or	6
rd	18
ti	0
uc	10
wo	16

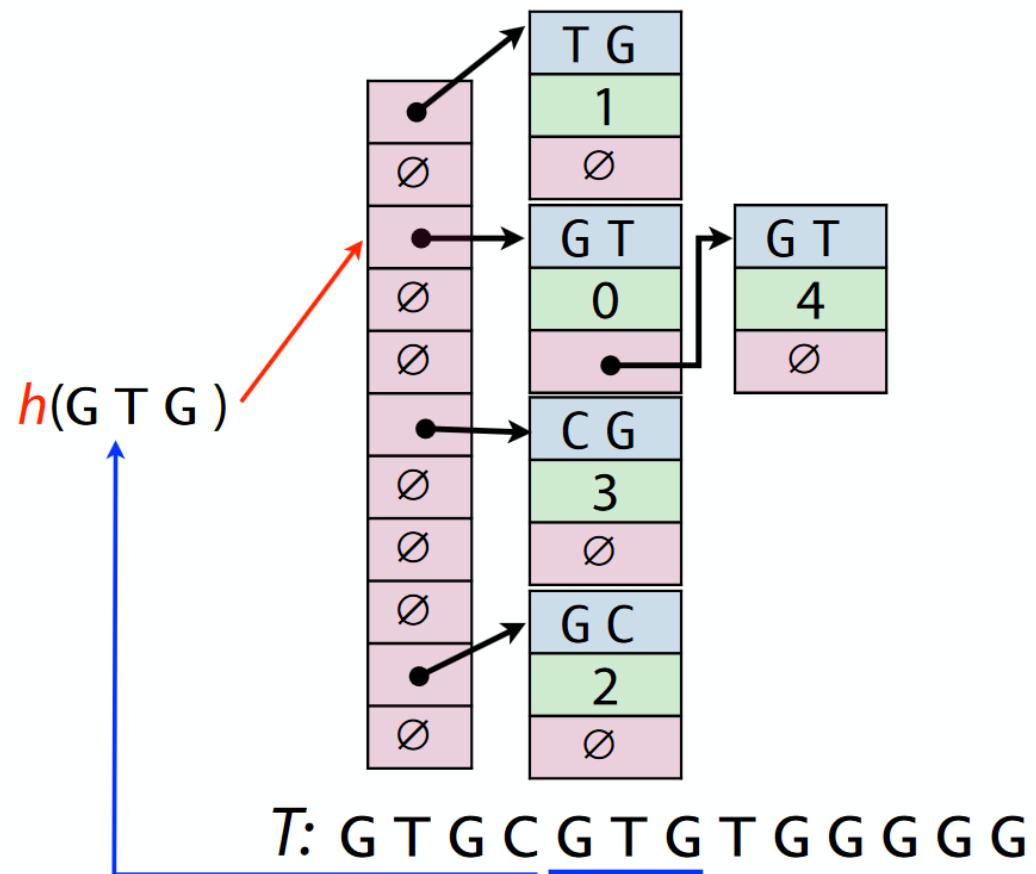
Trie:



Hash table:



Indexelés: Hash tábla áttekintés



Feladat

- Irjunk közelítő illesztést úgy, hogy az használjon indexeket! Adjuk vissza a nem egyező karakterek számát is!
- A házi feladat notebookjának linkjét küldjétek el e-mailben:
 - lkp@caesar.elte.hu
- A hiba, amibe órán ütköztünk a BoyerMoore approximáció során az volt, hogy a BoyerMoore előfeldolgozás nem működött 1 karakterre, ezt lekezeltem, itt a javított notebook:
 - https://notebooks.azure.com/lkpeter/libraries/bioinfcourse/html/bm_appr.ipynb