

# Karakterláncok és pontos illesztés

Lehotay-Kéry Péter

lepuaai at inf dot elte dot hu

Ben Langmead diasora alapján

([www.langmead-lab.org/teaching-materials](http://www.langmead-lab.org/teaching-materials))

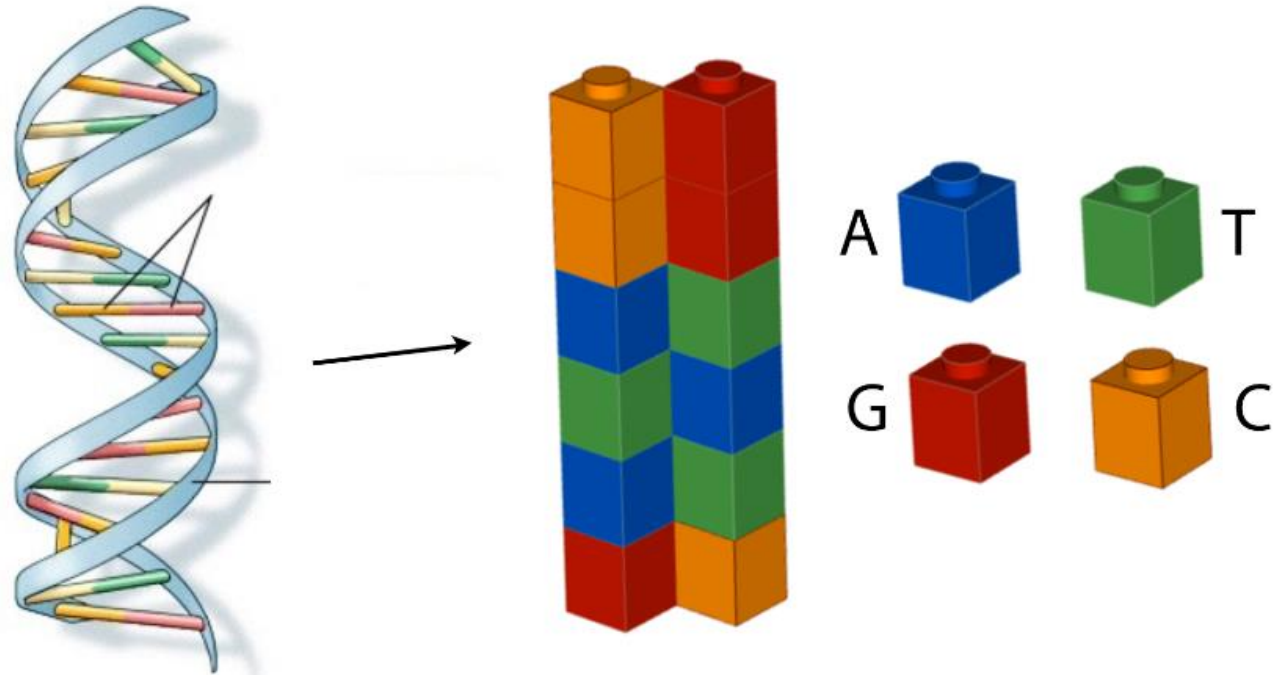
# Karakterlánc (String) definíciók

- $S$  string egy véges, rendezett karakterlista
- A karakterek egy  $\Sigma$  ábécéből kerülnek ki,  $|\Sigma|$  eleme van.
  - Nukleinsav ábécé:  $\{A, C, G, T\}$
  - Aminosav ábécé:  $\{A, R, N, D, C, E, Q, G, H, I, L, K, M, F, P, S, T, W, Y, V\}$
- $S$  hossza  $|S|$ , ahány karaktere van  $S$ -nek.
- $\epsilon$  az üres karakterlánc.  $|\epsilon|=0$

# Karakterlánc (String) definíciók

- $S$  és  $T$   $\Sigma$  ábécé fölötti karakterláncok, konkatenációjuk  $S$  karakterei, melyeket  $T$  karakterei követnek:  $ST$ .
- $S$  részstringje  $T$ -nek, ha léteznek  $u$  és  $v$  stringek, melyekre  $T=uSv$
- $S$   $T$  prefixe, ha létezik  $u$  string, melyre  $T=Su$ .
- $S$   $T$  szuffixe, ha létezik  $u$  string, melyre  $T=uS$ .
- $S$  valódi prefix (szuffix), ha se  $S$ , se  $u$  nem üres.
- Részszekvencia hasonló részstringhez, de a karakterek nem feltétlenül követik egymást: “cant” részszekvenciája “concatenate”-nek, de nem részstringje.

# Fordított komplement



U.S. National Library of Medicine

Double stranded  
DNA (double helix)

Double stranded  
DNA (lego version)

# Pontos illesztés

- Olyan helyeket keresünk, ahol  $P$  minta előfordul részstringként  $T$  szövegben. Minden ilyen hely egy előfordulás, vagy illeszkedés.
- Legyen  $n=|P|$  és  $m=|T|$  és  $n \leq m$ .
- Illesztés egy módja annak, hogy  $P$  karaktereit  $T$  karaktereivel szembe helyezzük. Ez lehet előfordulás vagy nem.

$P$ : word

$T$ : There would have been a time for such a word

Alignment 1: word

Alignment 2: word


# Pontos illesztés

*P*: word

*T*: There would have been a time for such a word

word word word word word word word word word word  
word word word word word word word word word  
word word word word word word word word word  
word word word word word word word word word  
word word word word word word word word word

One occurrence



*P*: word

*T*: There would have been a time for such a word

-----word-----word----->word  
----->----->----->

# Pontos illesztés: naive algoritmus

- Mennyi illesztés lehetséges adott  $n$ -re és  $m$ -re?
- Mennyi karakter összehasonlítás lehetséges maximum?
- Minimum?
- Mennyi karakterösszehasonlítás történik az alábbi példában?

$P$ : word

$T$ : There would have been a time for such a word

-----word----->word  
----->----->

$m - n$  mismatches, 6 matches

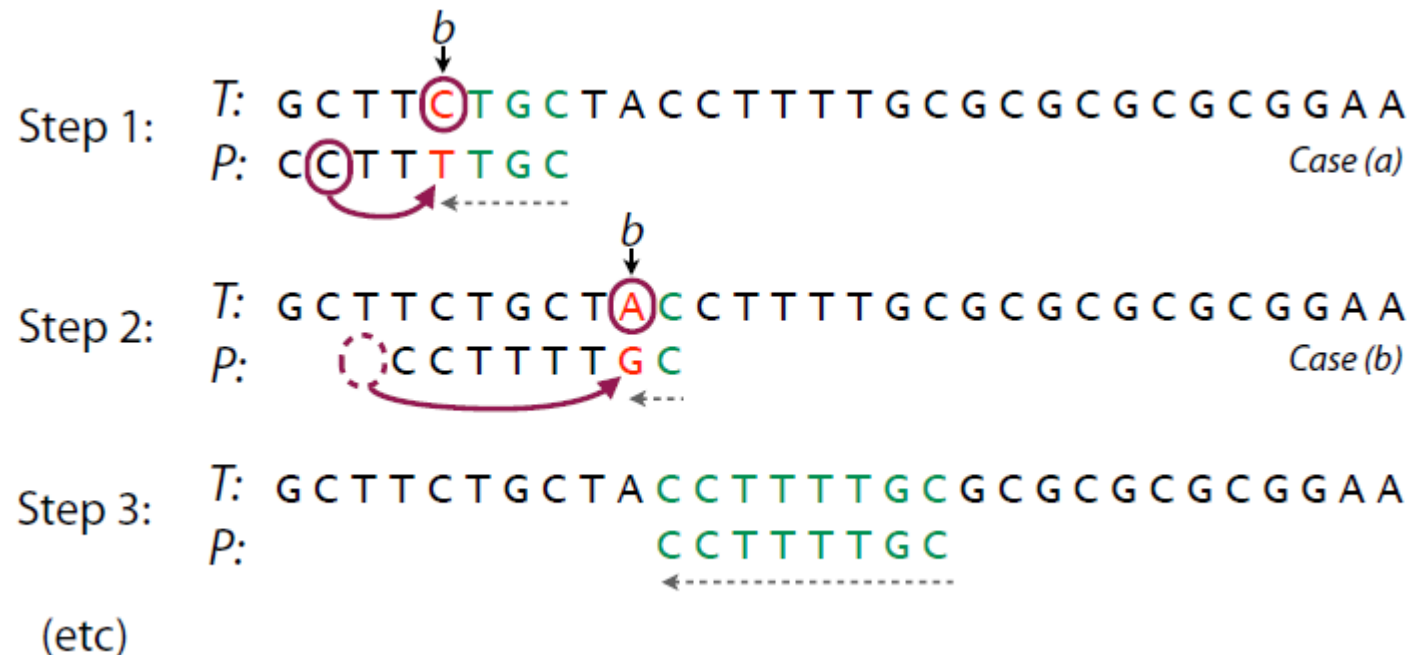
# Boyer-Moore

- Hagyjuk ki a következő illesztéseket, melyek biztosan nem fognak egyezni:
- Rossz karakter szabály: Ha nem egyezik a karakter, használjuk ki, hogy illesztéseket hagyjunk ki.
- Jó szuffix szabály: Ha egyezik néhány karakter, használjuk ki, hogy illesztéseket hagyjunk ki.
- Hosszabb kihagyásokért: Próbáljunk illesztéseket egy irányba, karakter összehasonlítást a másikba.



# Boyer-Moore: Rossz karakter szabály

- Ha nem egyezik egy karakter, legyen  $b$  a nem egyező karakter T-ben.  
Hagyjuk ki az illesztéseket, amíg
  - a)  $b$  nem egyezik a vele szembe levővel P-ben
  - B) P át nem halad  $b$ -n.



# Boyer-Moore: Jó szuffix szabály

- Legyen  $t$   $T$  részstringje, mely egyezett  $P$  szuffixével. Hagyjuk ki az illesztéseket, amíg
  - a)  $t$  nem egyezik a szemben levő  $P$  karakterekre
  - b)  $P$  prefixe nem egyezik  $t$  szuffixével

Step 1:  $T$ : CGTGCC **TAC** TTACTTACTTACTTACGCGAA  
 $P$ : CTTACTTAC *Case (a)*

Step 2:  $T$ : CGTGCC **TACTTAC** TTACTTACTTACTTACGCGAA  
 $P$ : CTTACTTAC *Case (b)*

Step 3:  $T$ : CGTGCCCTACTTACTTACTTACTTACTTACGCGAA  
 $P$ : CTTACTTAC

# Boyer-Moore: előfeldolgozás

- Előre kiszámított lépések a rossz karakter szabályra,  $P=TCGC$

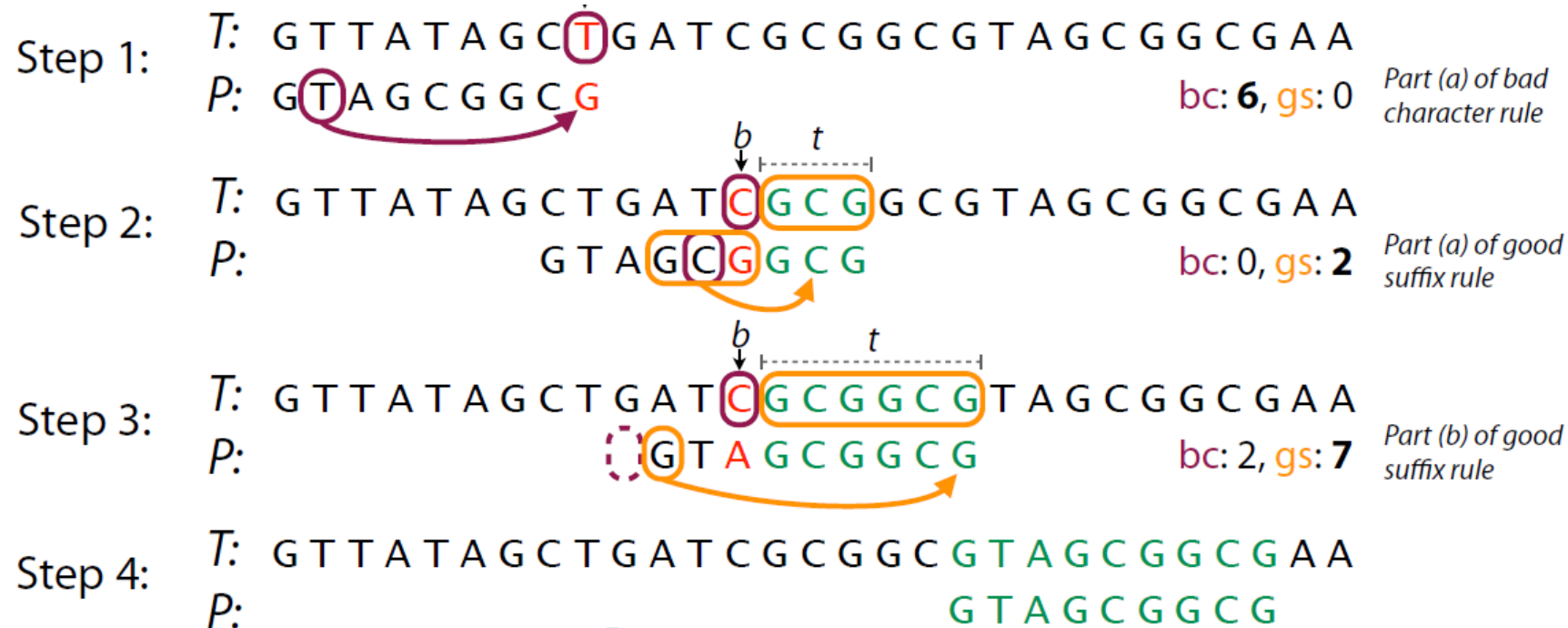
$P$

	T	C	G	C
A	0	1	0	3
C	0	-	2	-
G	0	1	-	0
T	-	0	1	2

$T$ : A A T C A A T A G C  
 $P$ : T C G C

# Boyer-Moore: Tegyük össze

- Minden illesztésre használjuk a rossz karakter, vagy jó szuffix szabályt, amelyek több illesztést hagy ki.



# Feladat

- Választható feladatok:
  1. Írjuk át a naive és Boyer-Moore illesztést úgy, hogy számolja össze, hány illesztést és hány karakter összehasonlítást végez! (2 szám lesz az eredmény az egyik lefutásakor és 2 szám a másik lefutásakor)
  2. Írjuk át a naive illesztést úgy, hogy a minta fordított komplementjeit is megtalálja! (Az eredmény az illeszkedő pozíciók és a fordított komplement illeszkedő pozíciói egy tömbben, szigorúan monoton növekvő sorrendben.)
- Készítsetek egy notebookot, mely a házi feladatokat fogja tartalmazni, ezt osszátok meg az e-mail címemmel: lepuaai at inf dot elte dot hu