

# Közelítő illesztés

Lehotay-Kéry Péter

lepuaai@inf.elte.hu

Ben Langmead diasora alapján

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

# Szükség van közelítő illesztésre

- Szekvencia eltérések miatt
  - Szekvenálási hibák
  - Természetes variációk

# Eltérés, behelyettesítés

*T*: GGAAAAGAGGTAGCGGCGTTTAAACAGTAG  
          | | | | | | | |  
*P*: GTACGGCG

# Beillesztés

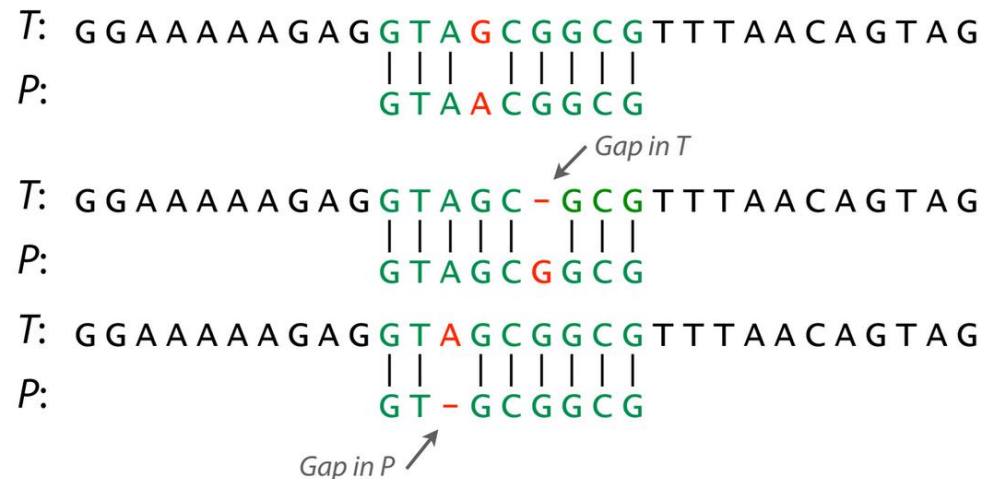
*T*: GGAAAAGAGGTAGC - GCGTTTAACAGTAG  
          | | | | | | |  
*P*: GTAGCGGCG

# Törlés

*T*: GGAAAAGAGGTAGCGGCGTTTAAACAGTAG  
          || - |||||  
*P*: GT - GCGGCG

# 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*

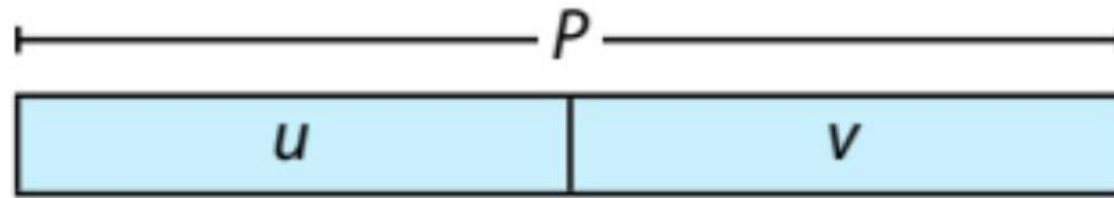
- Szerkesztés távolság a szükséges minimális számú szerkesztés, ahhoz, hogy egyik a másikká váljon.

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.



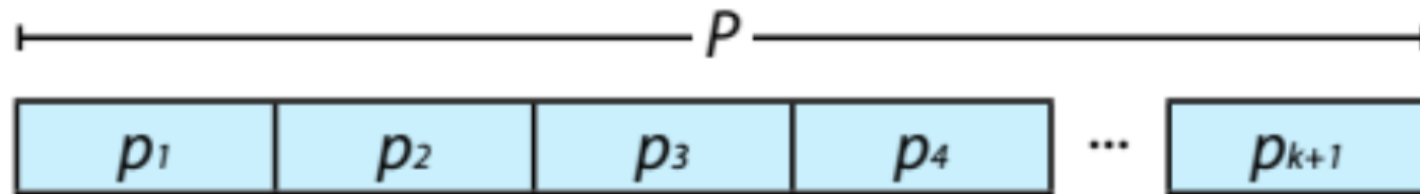
Either the edit goes here...

...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.