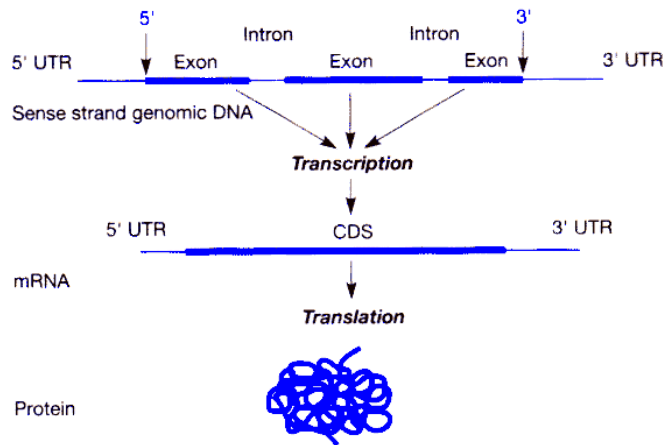


# 6. Nukleotidszekvenciák analízise

1. Bevezetés
2. DNS-szekvenciák összeszerelése
3. Génvadászati foratókönyv
4. A repetitív elemek kiszűrése
5. Hasonlósági keresések
6. ORF-detektálás
7. Kodonyakoriság-eltolódások detektálása
8. Funkcionális helyek detektálása
9. Integrált génelemzés
10. Webhelyek, példák

## Bevezetés

### Centrális dogma



Rövidítések: UTR: Untranslated Region, CDS: Coding sequence

### A genetikai kód

	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	
<b>T</b>	Phe	Ser	Tyr	Cys	T
	Leu		<b>Stop!</b>	<b>Stop!/Sec</b>	A
<b>C</b>	Leu	Pro	His	Arg	T
			Gln		A
	G	G			G
<b>A</b>	Ile	Thr	Asn	Ser	T
	<b>Met/Start</b>		Lys	Arg	C
<b>G</b>	Val	Ala	Asp	Gly	A
			Glu		G

A genetikai kód közel univerzális (apróbb eltérések azért vannak egyes mikroorganizmusokban, mitokondriumokban, stb.)

A TGA időnként nem Stop–ot, hanem szelenociszteint (Sec) kódol (21. aminosav!)

---

## DNS–szekvenciák összeszerelése

- cDNS szekvenálásánál 200–500 bázis hosszú darabokat kapunk (egy menetben ekkorát lehet szekvenálni, újabban azért max. 1000–et is)
- Szekvenálási hibák: kb. 5% (hibás bázisok; kimaradt/tévesen beszúrt bázisok: ún. fantom INDEL–ek)
- Emiatt mindkét szálát többször meg kell szekvenálni a megbízhatóság végett
- A teljes szekvenciát a darabokból kell összeszerelni és konszenzust konstruálni:

```
AACCGTTTACGAAACCAGGTGC
AACCGTTTACGAAACCAGGTGCGCGCCCGCGGAAT
AACCGTTTACGAAACCAGGTGC
(konszenzus: ) AACCGTTTACGAAACCAGGTGCGCGCGCGCGGAATCCTAAAAA
                                     CGCGCGCGCGGAATCCTAAAAA
                                     AGGGAATGCTGAGGAG
```

Kisbetűk: kisebb megbízhatóság

- Összeszerelés: különféle programokkal, pl. [TIGR Assembler](#)

---

## Génvadászati forgatókönyv

1. Kiszűrni a repetitív elemeket (kiderül, hol *nem* lehetnek a szabályozó és kódoló régiók)
2. Hasonlósági keresést végezni az ismert szekvenciák adatbázisain
3. Megvizsgálni a kodongyakoriságok eltolódásait (a fehérjekódoló régiók legegyértelműbb jele)
4. Funkcionális helyeket (iniciáció, termináció, exon–intron határok, promoterek, stb.) keresni az ismert szekvenciamintázatok alapján
5. A fenti lépésekből nyert információt egységes képbe foglalni, konszenzusos képet kialakítani

### Fontos:

- A legtöbb program specifikus valamely fajra vagy élőlénycsoportra
- A programokat a megfelelő kontextusban kell használni (pl. exon–intron határt ne keressünk cDNS–ben)

---

## A repetitív elemek kiszűrése

- Repetitív elemek: rövid (vagy hosszabb) ismétlődő szekvenciadarabok
- Eukariota genomok jelentős részét teszik ki
- "Interspersed" repeat: a genomban elszórva találhatóak, többnyire transzpozázibilis elemek (pl. transzpozonok) inaktiválódott kópiái
- Számos családjuk van, adatbázisuk: [RepBase](#)
- Minden DNS–elemzés előtt ki kell szűrni ezeket a szekvenciából, mert minden elemzőprogramot megzavarnának.
- Jelentősége: megmutatja, hol *nem* lehetnek szabályozó és kódoló régiók
- Programok erre: pl. [Censor](#), [RepeatMasker](#): a RepBase vagy hasonló adatbázis alapján szűrnek
- Pl. Censor futtatás eredménye:

**Bemenő szekvencia: humán kreatin kináz génjének részlete**

```
> HUMCKMM1
ggatccttcctcctctggcctcccaagtgtgggattacaggtgtgagccactgcactg
gcctattacccttctcaggctctggagtcacatccttctgctctgtctccctcagttcaat
tgtttttgtttttgttttttttagacacagctctcgctctgtcaccaaggctggagt
gcagcagtgcgatcacagctcaccgcagcctcacctcccaggctcaagtatcctccat
ctcggcctctgagtagctgagactataggtgtgtccacatgtcgggctaattttgtatt
tttagtagagacagggtttcaccgcgttggccaggggtgtctgaactcctgagctcaag
caatcctcctgctcagcctcctgtttttagatcccacaaataacttgtgatg
tttgcctttctatacctgggtcatttaacattttcttttctttctttcttttttttt
ttttttgtgagactgagctctgtctgtcactcaggctggagggaatggatcctcag
ctcactgcaacctccacctcctaggttcaagcaattcttatgcctcagcctcctggctag
ctgggat tacaggcgtgtgtcaccatgccaggctaattttttagttagtagatgg
ggtttaccatgttggccaggctggcttgaactcctggcctcaagtgtccaccgcct
ccgcctctgcctcccaaagtgtgggattacgggctgagccactgtgcccggccatc
aacattttcaactgtcaatcacaatgggattaaaactcctcccacagcccctagggacca
```



```

Forward 1:
      10      20      30      40      50
0 PLSLIPVTST LEDLVRGISD CESSLTERS Y !AGLVE!RDE KSDYCMQWLH
50 SVFLFQSIPE IPDDLKQLYK TVWEISQKTV LKM

Forward 2:
      10      20      30      40      50
0 H!ALYQ!HLH SKILSGEFQI VNPHELLKDLT ERGLWNEEMK NQIIACNGSI
50 QFSFFRAYQK FLMT!SNSIR PCGKSLRRLF SR

Reverse 0:
      10      20      30      40      50
0 HLENSLLRDF PHGLIELLQV IRNFWYALKK EN!MEPLHAI I!FFISSFHK
50 PRSVRSFSK! GFTI!NSPDK IFECRCYWYK AQW

Reverse 1:
      10      20      30      40      50
0 ILRTVF!EIS HTVL!SCFRS SGISGML!KR KTEWSHCMQ! SDFSSLHSTS
50 PAQ!DLSVSE DSQSEIPLTR SSSVDVTGIR LNG

Reverse 2:
      10      20      30      40      50
0 S!EQSSERFP TRSYRVASGH QEFLVCSEKG KLNGAIACNN LIFHLFIPQA
50 PLSKIFQ!VR IHNLKFP!QD LRV!MLLV!G SM

```

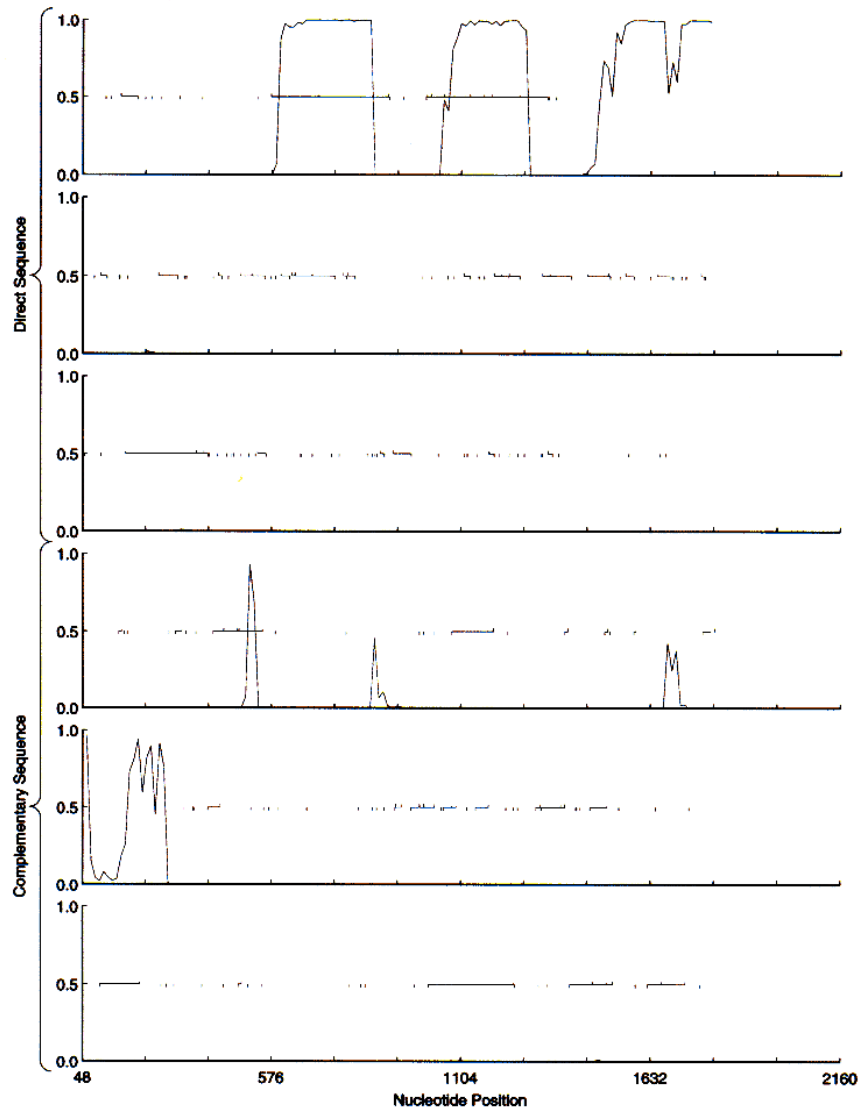
- A valódi leolvasási keret valószínűleg a leghosszabb megszakítatlan szekvenciát eredményez.
- A Start kodon azonosítása sokkal nehezebb (ATG a metionin kódja is), ld. később.

## Kodongyakoriság–eltolódások detektálása

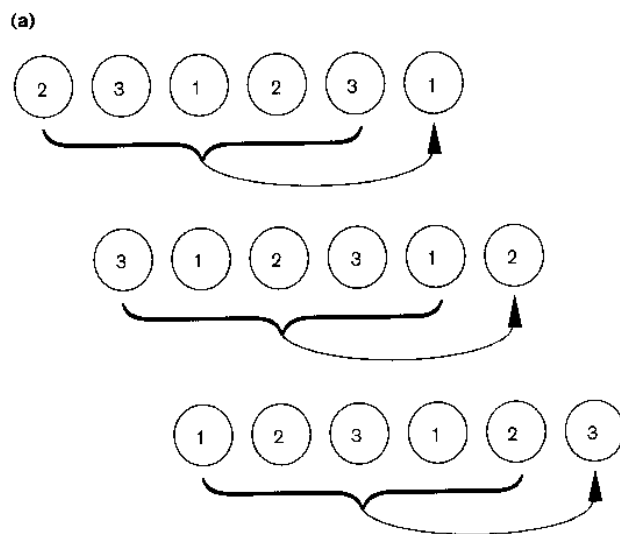
- A nemkódoló és a kódoló régiókban az egyes kodonok (nukleotidhármások) gyakorisága jellegzetesen eltérő
- Különböző organizmusok más–más kodonpreferenciát mutatnak. Pl. a szérin hatféle kodonjának gyakoriságai 5 fajban:

Kodon	<i>E. coli</i>	<i>D. melanogaster</i>	<i>H. sapiens</i>	<i>Z. mays</i>	<i>S. cerevisiae</i>
AGT	3	1	10	3	5
AGC	20	23	34	30	4
TCG	4	17	9	22	1
TCA	2	2	5	4	6
TCT	34	9	13	4	52
TCC	37	48	28	37	33

- Organizmuson belül is változhat! Pl. *E. coli*-nál három kategória
- Sokféleképpen meg lehet ragadni, legsikeresebb: a **hexamerek** (nukleotidhatások) előfordulási gyakorisága
- Pl. GeneMark program valószínűséget számol (kódoló–e a régió) a hexamergyakoriságok alapján, mind a 6 leolvasási keretben.



- Más programok (pl. GenScan) Markov–modellekkal dolgoznak:



Pozíciófüggő (inhomogén) ötödrendű Markov–modell: minden pozícióban valamely nukleotid valószínűségét az előtte levő öt nukleotid milyensége határozza meg, s ez a kerethez viszonyított pozíciótól is függ. A valószínűségértékeket a már ismert gének alapján határozzák meg, új szekvencia valószínűsége a modell alapján kiértékelhető.

# Funkcionális helyek detektálása

Promoterek, exon–intron határok, translációiniciációs helyek, terminációs helyek azonosítása, az ezekre jellemző szekvenciamintázatok alapján

## Az algoritmusok jóságának két paramétere:

- **Érzékenység:** a valóságos funkcionális helyek mekkora részét jósolja meg helyesen
- **Specifititás:** a megjósolt helyek mekkora része valóságos

## Különböző megközelítések:

- Konszenzusszekvencia alapján (a leggyakoribb nukleotidok megadása): nem megbízható a variabilitás miatt
- Súlymátrix megadása: minden pozícióra közöljük a 4 nukleotid gyakoriságát, ehhez hasonlítjuk az új szekvenciát. Korlátozott mértékben sikeres.
- Bonyolultabb matematikai modellek (pl. rejtett Markov–modellek). Ígéretes, de nem kiforrott.

## Promoterek

- Számos program, az ismert promoterek szekvenciáinak könyvtára alapján
- Pl. TATA–box, "cap" szignál, stb.
- Kevésbé sikeres, mert pl. az emberi gének 30%–ánál nem találhatóak meg a szokásos promoterszignálok
- A sikeresség nem nagyobb, mint a kodongyakoriság–eltolódások alapján

## Poliadenilációs szignál

- A gének végén, ált. AATAAA konszenzusszekvenciával, melyet egy bonyolultabb szignál követ
- A gének felénél hiányzik
- Nem megbízható

## Transzlációs szignálok: start és stop kodon

- ATG (RNS–ben AUG) a start kodon, de a metionin kódja is
- **Kozak–féle szabályok:** a start kodon sokszor jellegzetes környezetben van, pl.

CCGCCAUGG

és hasonló (súlymátrixszal leírható)

- Gyenge az érzékenység (70–80%), de nem genomiális, hanem cDNS esetében jól használható
- Stop kodon: egyértelmű

## Exon–intron határok

- Az exon–intron határokon (5' vég a donor hely, 3' vég az akceptor) jellegzetes szekvenciák vannak, melyeket a spliceoszóma felismer
- Intronok szinte mindig GT–vel kezdődnek és AG–vel végződnek, ez kevés, de segít leszűkíteni a lehetőségeket
- Bonyolultabb mintázatok súlymátrixa alapján elég hatásos detektálás, de nem tökéletes
- Bonyolultabb matematikai modellek (pl. rejtett Markov) jelentős, de nem drámai javulást eredményeznek
- Leghatásosabb: együttesen alkalmazni a kodongyakoriság–vizsgálattal (80% fölé emelkedik az érzékenység)

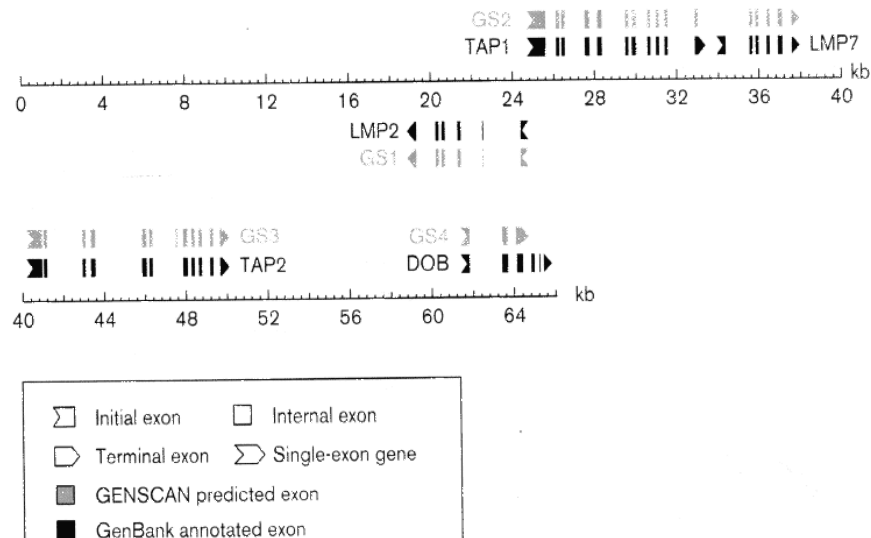
---

## Integrált génelemzés

- Az egyes predikciók hatékonyságát nagyon jelentősen növeli, ha a különböző algoritmusokból származó információt együttesen vesszük figyelembe. Ez vezetett az integrált génelemző programok kifejlesztéséhez
- A fenti eljárásokat egymás után alkalmazzák
- Számos ilyen program, szerver van, lásd pl. [Gene Prediction Services](#), [Gene Recognition Programs](#)
- 90% körüli érzékenységet, ill. specifitást lehet elérni velük.

## Korlátaik

- Csak néhány organizmusra létezik ilyen program
- Néhány kivétellel csak akkor működnek jól, ha csak egyetlen gén van a beadott szekvenciában
- Legtöbbször nagyon érzékenyek a szekvenálási hibákra
- Nem képesek kezelni az olyan eseteket, mint pl. alternatív splicing, átfedő gének, különleges promoterek, stb.
- Ajánlott *mindegyik* programot használni (más–más algoritmusokkal dolgoznak)
- Pl. a Genscan eredményessége egy 66 kilobázis méretű szakaszon:



Fekete jelöli a valóságot, szürke a GenScan általi jóslást (ld. jelmagyarázat). Az LMP2 gén jóslása tökéletes, de a TAP1–et és az LMP7–et a GenScan egybeolvasztotta, a TAP2–be eggyel több exont jóslott, a DOB–nak pedig túl korán véget vetett. Az exonok túlnyomó részét mégis hibátlanul detektálta.

---

## Webhelyek

Ld. pl. [Gene Prediction Services](#), [Gene Recognition Programs](#)

---