

Metódy v bioinformatike, 1-BIN-301/2-AIN-501

Vyučujú:

Broňa Brejová, M-163, brejova@fmph.uniba.sk

Tomáš Vinař, M-163, vinar@fmph.uniba.sk

Web: <http://compbio.fmph.uniba.sk/vyuka/mbi/>

Diskusné fórum a oznamy: Facebook

Literatúra:

I-INF-D-23: Durbin, Eddy, Krogh, Mitchison: Biological sequence analysis. Cambridge University Press 1998.

I-INF-Z-2: Zvelebil, Baum: Understanding Bioinformatics. Taylor&Francis 2008.

Skriptá k predmetu a poznámky na webstránke.

Ciele predmetu

- **Všetci:** Prehľad základných metód na výpočtovú analýzu biologických sekvencií a ďalších dát v molekulárnej biológii.
- **Informatici:** Algoritmy a dátové štruktúry, strojové učenie, pravdepodobnosť. Ako prejsť od problému v reálnom svete k matematickej abstrakcii.
- **Biológovia:** Matematické modely tvoriace základ populárnych bioinformatických nástrojov, používanie nástrojov, interpretácia výsledkov.
- **Všetci:** Skúsenosť s interdisciplinárnou spoluprácou.

Známkovanie

3 domáce úlohy 30% (10% každá)

Journal club 10%

Skúška alebo projekt 60%

Hodnotenie: A: 90+, B: 80+, C: 70+, D: 60+, E: 50+

- Dve verzie otázok: biologická a informatická
- Journal club: čítanie 1 článku v skupine a správa (prípadne nepovinná prezentácia)
- Projekt povinný pre doktorandov, nepovinný pre ostatných
- Na skúške povolený ťahák 2 listy A4
- Neopisovať!

Časy a miestnosti

Prednáška štvrtok 15:40-17:10 F1-328

Cvičenia informatici štvrtok 14:00-15:30 F1-328

Cvičenia biológovia štvrtok 17:20-18:50 F1-328 a M-217
(počítačová učebňa)

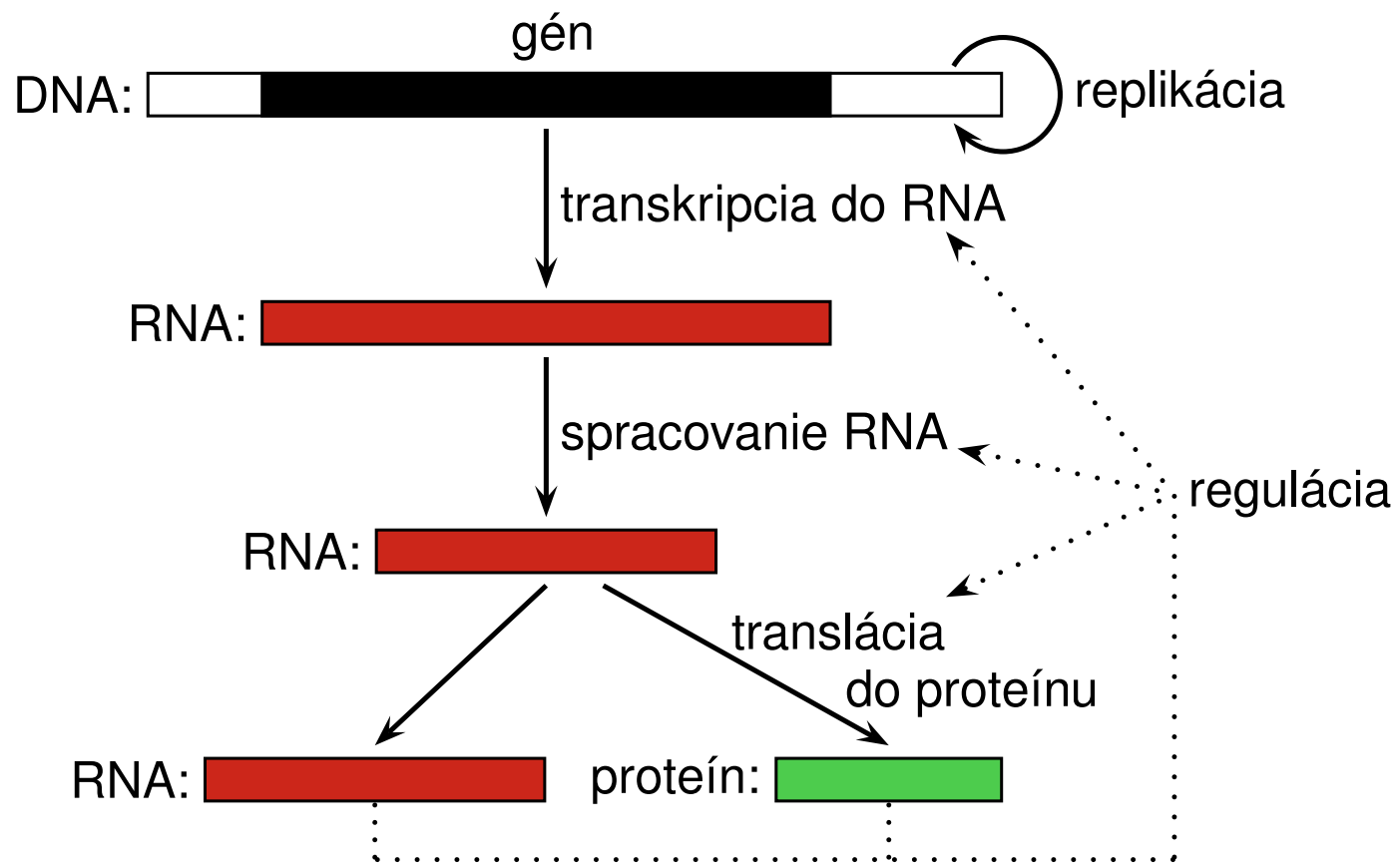
Čo nás v tomto predmete čaká

Typická prednáška

- Biologické pozadie problému
- Formulácia ako informatický problém
- Idea algoritmu (riešenia problému)

Typické cvičenia

- Informatici: ďalšie detaily algoritmov, potrebné poznatky z biológie
- Biológovia: aplikácia na konkrétne dáta, význam rôznych parametrov, potrebné poznatky z informatiky



- **Sekvenovanie, zostavovanie sekvencií**

Ako sa získavajú DNA sekvencie, akú rolu v tom hraje informatika?

Grafové algoritmy, skladačka z obrovského množstva malých kúskov

- **Zarovnávanie (sequence alignment)**

Ktoré sekvencie sú podobné na moju obľúbenú sekvenciu?

Čo presne robí BLAST a ako mu nastavím parametre?

Dynamické programovanie, heuristiky a ako povedať niečo o tom, ako dobre fungujú

- **Hľadanie génov**

Koľko génov má človek?

Skryté Markovove modely (HMM)

- **Evolúcia, rekonštrukcia fylogenetických stromov**

Ku komu máme bližšie: ku psom alebo k myšiam?

A má hroch bližšie k veľrybám alebo k prasatám?

Pravdepodobnostné modely evolúcie, princíp úspornosti (parsimony), metódy riešenia ťažkých úloh

- **Komparatívna genomika**

Ako sa líšime od šimpanzov?

Ktoré časti genómu sa vyvíjajú pomalšie alebo rýchlejšie ako by sme čakali a prečo?

Spojenie HMM a evolučných modelov

- **Expresia a regulácia génov**

Ktoré gény slúžia ako iniciátori bunkovej samodeštrukcie?

Dá sa jednoduchým vyšetrením rozlíšiť, či má konkrétny pacient zhubnú rakovinu a či konkrétny liek bude fungovať?

Zhlukovanie (clustering), biologické siete a ich vlastnosti

- **Transkripčné faktory**

Ako funguje mechanizmus riadenia expresie génov?

Vieme niektoré gény umelo zapínať a vypínať?

Hľadanie opakujúcich sa motívov v sekvenciách

Rozpoznávanie známych motívov pomocou strojového učenia

- **Štruktúra a funkcia proteínov**

Akú 3D štruktúru má môj obľúbený proteín?

Akú funkciu má v živej bunke?

Čo spôsobuje Alzheimerovu chorobu a prečo?

HMM, energetické modely, molekulárne dynamické simulácie

- **RNA**

Ako predikovať sekundárnu štruktúru RNA a hľadať RNA gény?

Koľko tRNA je v mojom obľúbenom genóme?

Dynamické programovanie, stochastické bezkontextové gramatiky

- **Populačná genetika**

Prečo Tibeťania nemajú problémy so životom vo veľkých výškach?

Ako to, že plemená psov vyzerajú tak rôzne, a napriek tomu sú jeden druh? A odkedy je vlastne pes najlepším priateľom človeka?

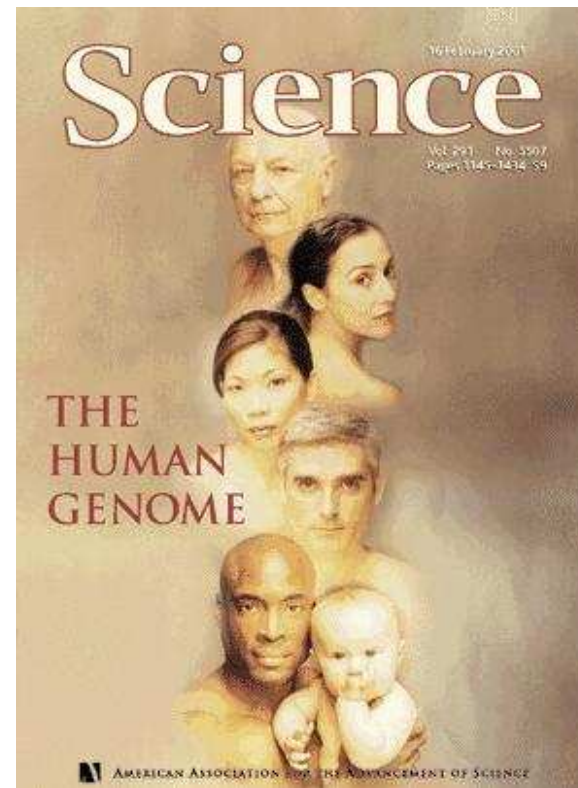
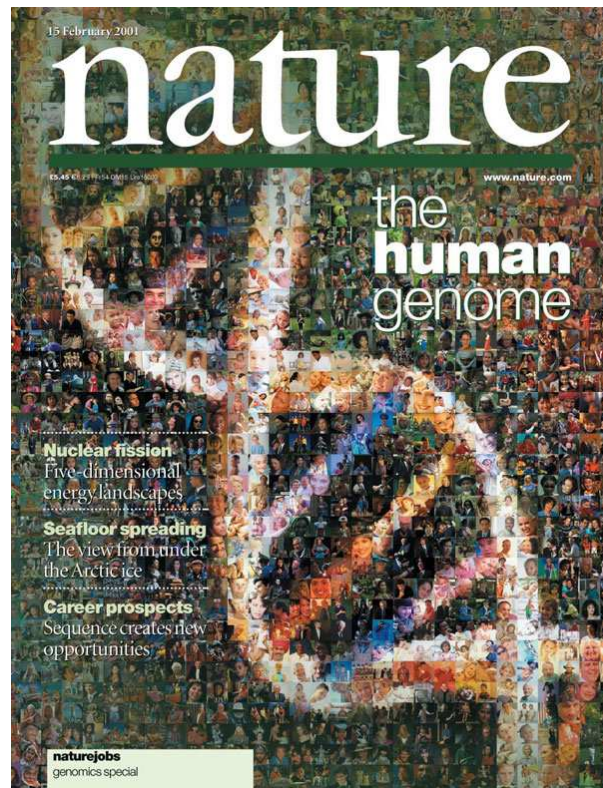
Pravdepodobnostné modely, stochastické algoritmy, štatistika

Budúci týždeň

- Nebude prednáška
- Cvičenia pre informatikov začnú normálne o 14:00
- Cvičenia pre biológov začnú o 15:40 v F1-328 (v čase prednášky), budú asi trvať o niečo viac ako 90 minút

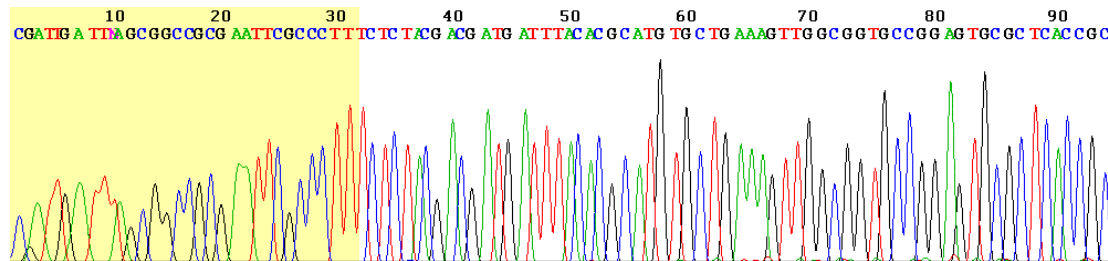
Sekvenovanie a zostavovanie genómov (genome sequencing and assembly)

Tomáš Vinar̃
22.9.2016



Sangerovo sekvenovanie

- Výsledok: sekvenovací profil (trace)



- Ďalej sa spracuje pomocou programu PHRED:
 - Na každej pozícii (kde sa dá) určí bázu (A,C,G,T)
 - Pre každú bázu odhadne kvalitu q
($10^{-q/10}$ je pravdepodobnosť chyby,
t.j. bázy s kvalitou $q > 40$ sú správne na 99.99%)
- Sangerovo sekvenovanie produkuje čítania (reads)
dlhé 500-1000 bp
- Ako osekvenovať dlhú DNA sekvenciu?

Typický priebeh sekvenovania

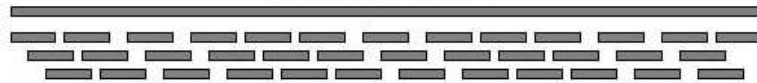
1. Chromozómy náhodne rozsekáme na menšie kúsky
(napr. pomocou sonikácie)
2. Menšie kúsky namnožíme
(napr. pomocou PCR, bakteriálneho klonovania a pod.)
3. Konce týchto kúskov osekvenujeme niektorou zo sekvenovacích technológií \Rightarrow mnoho krátkych reťazcov, ktoré nazývame **čítania**
4. Čítania **výpočtovo zostavíme** späť do chromozómov

Prehľad sekvenovacích technológií

Technológia	Dĺžka čítania	Chybovosť	Za deň	Cena za MB
1. generácia				
Sanger	do 1000 bp	< 2%	3 MB	\$4000
2. (next) generácia (cca 2004)				
Illumina	2× 150bp	< 2%	30 GB	< \$1
3. generácia (práve nabieha)				
PacBio	6-25 Kbp	15%	1 GB	\$2
Oxford Nanopore	up to 100kbp	30%		

Bioinformatický problém: Zostavenie genómu (sequence assembly)

- **Vstup:** krátke čítania sekvenovanej DNA
- **Cieľ:** zostaviť pôvodnú DNA
 - riadime sa zhodou v prekrývajúcich častiach čítaní
- Dôležité faktory:
 - **dĺžka genómu**
 - **pokrytie** (coverage) – koľko krát čítania pokrývajú genóm?



Formulácia problému

Najkratšie spoločné nadslovo (shortest common superstring)

Úloha: Daných je niekoľko reťazcov (čítaní), nájdite **najkratší** reťazec, ktorý obsahuje **všetky** vstupné reťazce ako (súvislé) **podreťazce**.

Motivácia: čo najviac využiť prekryvy medzi čítaniami

Príklad:

Vstup: GCCAAC, CCTGCC, ACCTTC

Výstup: CCTGCCAACCTTC (najkratšie možné)

Najkratšie spoločné nadslovo

- **Problém je NP ťažký**
takže nepoznáme rýchly algoritmus, ktorý vždy nájde najlepšie riešenie
- **Jednoduchá heuristika:** opakovane nájdí dva reťazce, ktoré sa prekrývajú najviac a zlúč ich do jedného reťazca
- Príklad: CATATAT, TATATA, ATATATC
Optimum: CATATATATC, dĺžka 10
Heuristika: CATATATCTATATA, dĺžka 14
- V skutočnosti táto heuristika **aproximačný algoritmus:**
Nájdene riešenie je najviac $3,5\times$ horšie ako optimálne
T.j. je to 3,5-aproximačný algoritmus
(možno aj 2-aproximačný, otvorený problém)
- Existuje aj 2,5-aproximačný algoritmus

Najkratsie spoločné nadslovo: Čo sme nezahrnuli do formulácie

- V sekvenovaní sa vyskytujú chyby (cca 1 zo 100 báz)
- Polymorfizmus
- Orientácia čítaní (vlákno, strand)
- Kontaminácia cudzou sekvenciou (napr. baktérie, v ktorých sa segmenty DNA klonovali), chiméry
- Viac chromozómov, neúplné pokrytie čítaniami
- Repetitívna sekvencia (sequence repeats, opakovania)

cca 50% ľudského genómu

Príklad: 10xTTAATA, 10xATATTA, 3xTTAGCT

TTAATATTAGCT?

TTAATATTAATATTAATATTAATATTAGCT?

TTAATATTA + ATATTAGCT?

Najkratšie spoločné nadslovo: Zľahčujúce faktory

Prídavná informácia: spárované čítania (Pair-end reads)



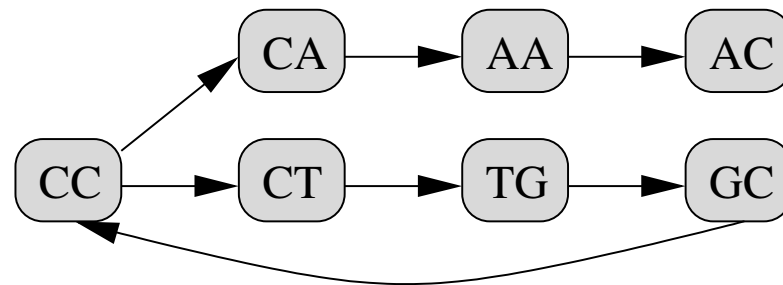
Zjednodušenie: nemusíme spojiť všetko do jedného reťazca, spájame len časti spojené viacerými čítaniami \Rightarrow konzervatívny prístup (radšej menej pospájať, ale nerobiť chyby)

Najkratšie spoločné nadslovo: Zhrnutie

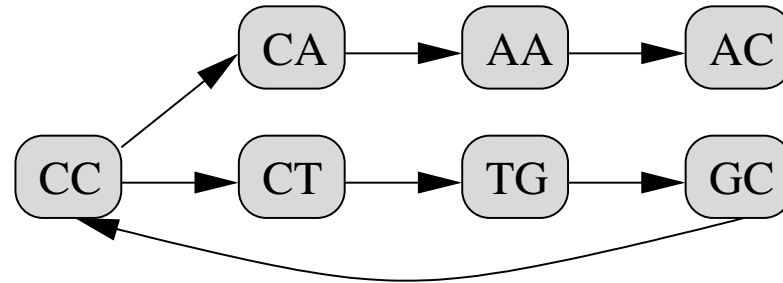
- Nerealistická formulácia, ťažký výpočtový problém
- Ale teoretický problém môže poskytnúť nejaký posun k pochopeniu skutočného problému
- Overlap-Layout-Consensus prístup motivovaný greedy algoritmami pre najkratšie spoločné nadslovo

Skladanie krátkych čítaní: de Bruijnov grafy

- Predpokladajme jednu orientáciu, žiadne chyby, jeden chromozóm úplne pokrytý čítaniami
- Nasekajme čítania na (prekrývajúce sa) kúsky dĺžky k
- Zostavme z nich **de Bruijnov graf**
 - **vrcholy**: podreťazce dĺžky k všetkých čítaní
 - **hrany**: nadväzujúce k -tice v rámci každého čítania (s prekryvom $k - 1$)
 - Graf je orientovaný (hrany majú smer)
- **Príklad**: $k = 2$, čítania: CCTGCC, GCCAAC



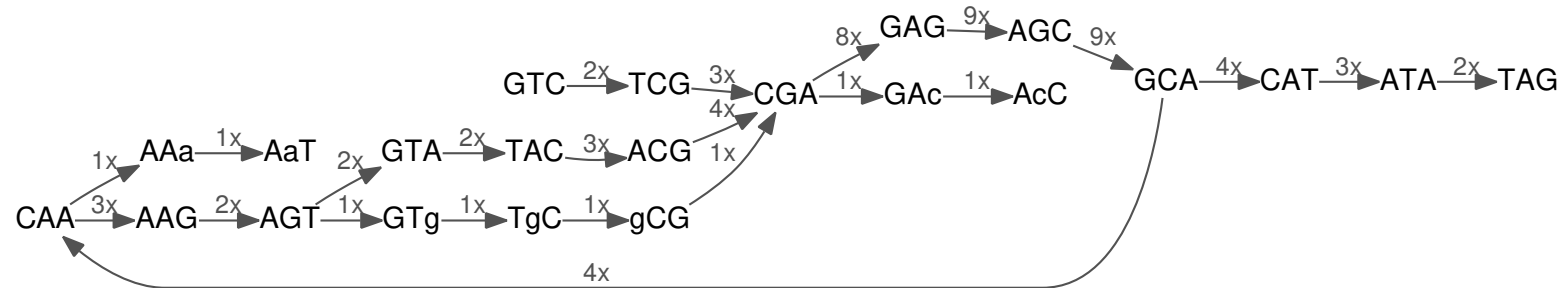
Ako použiť de Bruijnovu grafy?



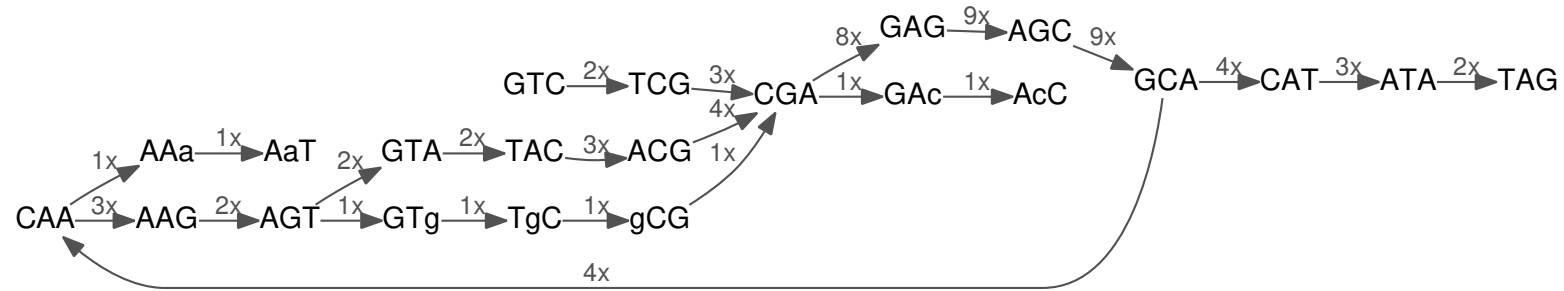
- jediný chromozóm a žiadne “nejednoznačné” k -tice
⇒ zostavenie = **Eulerovská cesta**
(cesta v grafe, ktorá použije každú hranu práve raz)
- Eulerovskú cestu možno nájsť v čase $O(m + n)$
- v realistickom prípade:
zostavenie genómu zodpovedá niekoľkým
pochôdkam v de Bruijnovom (nazývame **kontigy**),
ktoré dohromady pokrývajú veľkú časť hrán

Príklad: sada čítaní a zodpovedajúci deBruijnov graf

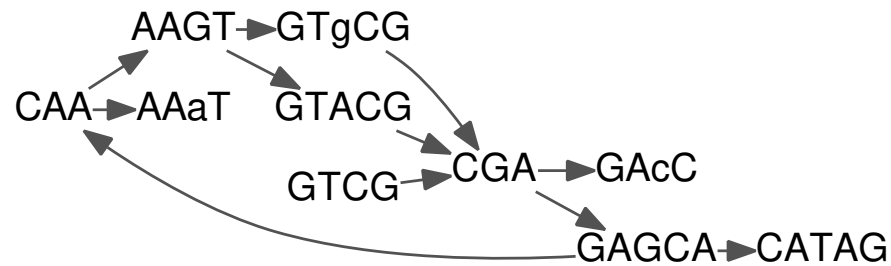
GTCGAGCAAGTACGAGCATAG
 TCGAGCA AGCATAG
 AGCAAaT AGCATAG
 GTCGAcC GTACGAG
 GTCGAGC TACGAGC
 CGAGCAA ACGAGCA
 AGTgCGA
 CAAGTAC
 GCAAGTA GAGCAT
 GAGCAAG GAGCATA
 TACGAGC



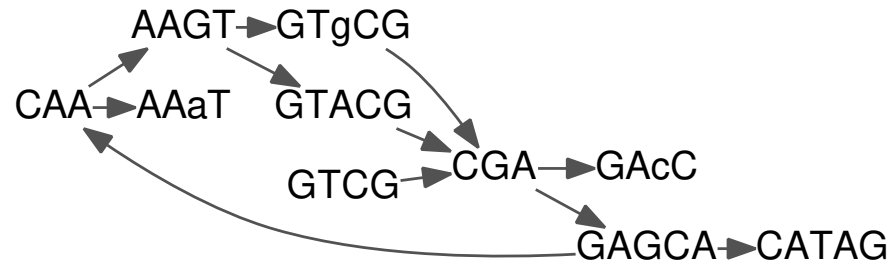
Príklad: zjednodušovanie de Bruijnovho grafu



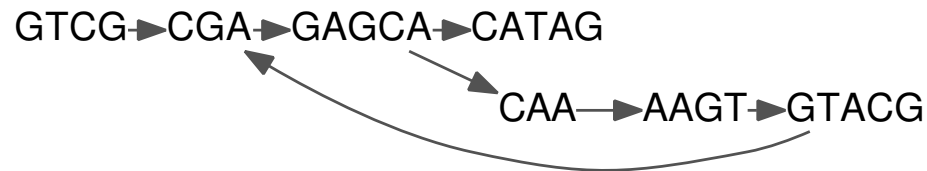
Spojíme jednoznačné cesty do vrcholov



Príklad: odstraňovanie chýb z de Bruijnovho grafu



Odstránenie chýb (výbežkov a bublín s nízkym pokrytím)



Spájaním dostaneme 4 kontigy (pôv. GTCGAGCAAGTACGAGCATAG)



Typické výsledky zostavovania genómov

- Veľa **kratších kontigov**, ktoré možno ďalej kombinovať do väčších celkov (**scaffolds**) pomocou ďalšej informácie (napr. spárované čítania)
- Niektoré časti nemožno jednoznačne zostaviť z dôvodu **dlhých opakujúcich sa sekvencií**

Príklad: človek chr14, 88 Mbp, 70× pokrytie

Metóda	Počet kontigov	Chýb	N50 po korekcii
Velvet (základný de Bruijn)	>45000	4910	2.1 kbp
Velvet (scaffolding)	3565	9156	27 kbp
AllPaths-LG	225	45	4.7 Mbp

N50: kontigy s touto alebo dlhšou dĺžkou pokrývajú 50% genómu

korekcia: rozsekne všetky zle spojené kontigy

História sekvenovania genómov

- 1976 MS2 (RNA vírus) 40 kB
- 1988 projekt sekvenovania ľudského genómu (15 rokov)
- 1995 baktéria H. influenzae 2 MB, shotgun (TIGR)
- 1996 S. cerevisiae 10 MB, BAC-by-BAC (Belgicko, Británia)
- 1998 C. elegans 100 MB, BAC-by-BAC (Wellcome Trust)
- 1998 Celera: ľudský genóm do troch rokov!
- 2000 D. melanogaster 180 MB, shotgun (Celera, Berkeley)
- 2001 2x ľudský genóm 3 GB (NIH, Celera)
- po 2001 Myš, potkan, kura, šimpanz, pes, makak,...
- 2007 Watsonov a Venterov genóm (454)
- 2012 1000 ľudských genómov
- čoskoro 10k genómov stavovcov

Zhrnutie

- Sekvenovanie genómu je zložitý proces, v ktorom hrá bioinformatika dôležitú úlohu
- V súčasnosti niekoľko nových technológií, nízka cena, krátke čítania
- Problém zostavovania genómu, najkratšie spoločné nadslovo
- Overlap-Layout-Consensus
- Praktické riešenie: de Bruijnove grafy
- V zostavenej sekvencii môžu byť chyby, medzery, viaceré kontigy
- Pokrytie genómu a veľkosť čítania hrajú najdôležitejšiu úlohu pri tom, ako fragmentovaný bude výsledok:
 - pre Sanger: 7-10× pokrytie
 - pre NGS: 40-70× pokrytie

Použitie NGS: Populačná genetika

- Sekvenujeme krátke čítania z genómu určitého človeka
- Ako sa môj vlastný genóm líši od genómu “priemerného” človeka?
- Ako jednoduché genetické rozdiely ovplyvňujú fenotyp?
- Personalizovaná medicína
- Populačná štruktúra, história ľudstva
- Etické otázky

Problémy:

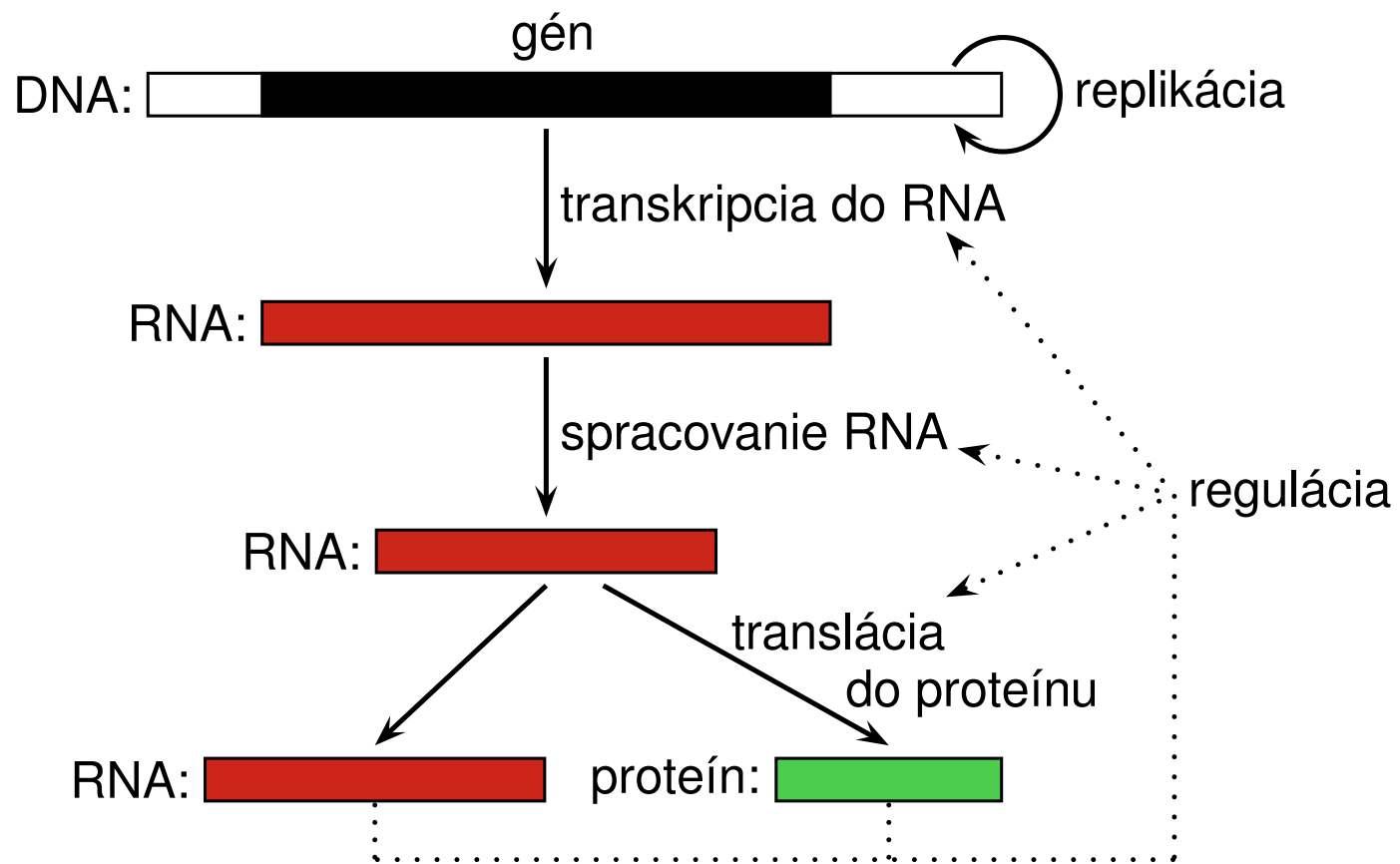
- Mapovanie krátkych čítaní na referenčnú sekvenciu
- Identifikácia rozdielov (malých a väčších)

Použitie NGS: Environmentálne sekvenovanie – Metagenomika

- Aké mikroorganizmy žijú v našich telách?
črevná a žalúdočná flóra, ústna dutina, koža, ...
- Diverzita mikroorganizmov v rôznych ekosystémoch
- Ťažké izolovať jednotlivé organizmy
- Sekvenujeme zmes čítania z rôznych genómov
- Snažíme sa zostaviť aspoň krátke kontigy

Problémy:

- Oddelenie čítaní/kontigov patriacich do rôznych genómov



Použitie NGS: Hľadanie génov, väzobných miest,...

- Sekvenovať môžeme aj RNA, dostávame gény v genóme
- Chip-Seq: vyfiltrujeme kúsky DNA, na ktoré je naviazaný určitý proteín, sekvenujeme, mapujeme na genóm

Problémy:

- Identifikácia miest zostrihu
- Identifikácia väzobných miest podľa hĺbky pokrytia

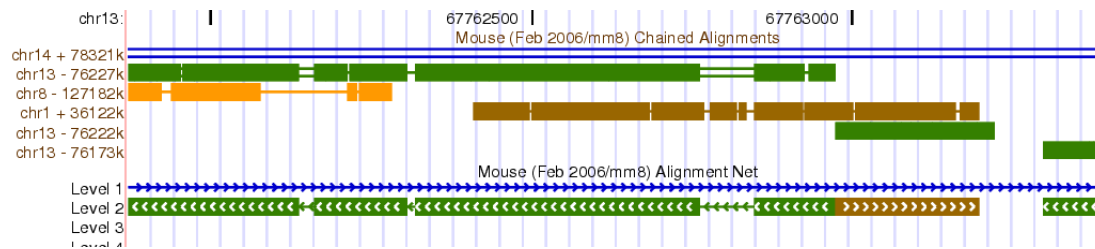
Oznamy

- Výber článku na journal club formulárom na stránke do stredy 19.10.
- Prvá domáca úloha: zadanie na budúci týždeň, čas na vypracovanie cca 2 týždne (detaily nabudúce)
- Ďalšie info <http://compbio.fmph.uniba.sk/vyuka/mbi/>

Zarovňávanie sekvencií (sequence alignment) 1/2

Tomáš Vinar

6.10.2016



[Durbin et al., 1998, kapitola 2]

Problém: Lokálne zarovnanie (local alignment)

ggcccttggagttgactgtcctgctgctccttgagg
ccattctcagagagaggaagtggcctcattttaatc
cgcttcccacagccttgtcctttccagacccatggg
agagggaggggctgagggtgtggctgagcccacca
agtcacgcgctcactctgcaggtccctctccccaag
gccgtggccttgggagcccgtggatcccagtgagtg
acgcctccacccccgacctactcgggcagtttaac
ccttgttgttcacttgcagacatcgtgaacacggcc
cggcccagcagagaaggccataatgacctatgtgtcc
agcttctaccatgccttttcaggagcgcagaaggta
ccgagcagggccaggcaggccctcctcgccgccacc
gcgcaatgccgcccgtgcctctcgctcccgtgctc
acctcatttctcttgcagacggcagtggcctctctc
caactggaagccacccccagctccct...

tgatgccgaggatgtgttcgctcgagcatccggacga
gaagtcacacctacgtggtcacactactatcacta
cttagcaaaactcaagcaggagacggtgcagggcat
aagcgtatcggttaagggtggcattgccatggag
aacgacaaaatgggtccacgactacgagaacttcaca
agcgatctgctcaagtggatcgaaacgacccatccag
tcgctgggagcagcgggagttcgaaaactcgctggcc
ggcgtccaagggcagttggcccagttctccaactac
cgcacccatcgagaagccgcccagtttgtggaaaag
ggcaacctcgaggtgctccttttcacctgcagtc
aagatgcgggccaacaaccagaagccctacacacc
aaagagggcaagatgatttcggacatcaacaaggcc
tgggagcgtctggagaaggccgagcacgaacgcgaa
ttggccctgcgcgaggagctcatccg...

Vstup: dve sekvencie

Problém: Lokálne zarovnávanie (local alignment)

ggccttggagttgactgtcctgctgctccttgagg
ccattctcagagagaggaagtggcctcattttaatc
cgcttcccacagccttgtcctttccagacccatggg
agagggaggggctgagggtgtggctgagcccacca
agtcacgcgtcactctgcaggtccctctcccccaag
gccgtggccttgggagcccgtggatcccagtgagtg
acgcctccacccccgcctactcgggcagtttaac
ccttgttgttcacttgcagacatcgtgaacacggcc
cggcccgacgagaaggccataatgacctatgtgtcc
agcttctaccatgccttttcaggagcgcagaaggta
ccgagcagggccaggcaggccctcctcgccgccacc
gcgcaatgccgccgctgcctctgcctcccgtgctc
acctcatttctcttgacagcggcagtggcctctctc
caactggaagccacccccagctccct...

tgatgccgaggatgtgttcgctcgagcatcggacga
gaagtccatcacctacgtggtcacctactatcacta
cttagcaaaactcaagcaggagacggcgcagggcat
aagcgtatcggtaagggtggcggcattgccatggag
aacgacaaaatgggtccacgactacgagaacttcaca
agcgatctgctcaagtggatcgaaacgaccatccag
tcgctgggagcgggagttcgaaaactcgctggcc
ggcgtccaagggcagttggcccagttctccaactac
cgcaccatcgagaagccgccaagtttgtggaaaag
ggcaacctcgaggtgctccttttcacctgcagtcc
aagatgcgggccaacaaccagaagccctacacacc
aaagagggcaagatgatttcggacatcaacaaggcc
tgggagcgtctggagaaggccgagcacgaacgcgaa
ttggcctgcgcgaggagctcatccg...

Výstup: podobné úseky (zarovnanie, alignments).

```
CCCGACGAGAAGGCCATAATGACCTATGTGTCCAGCTTCTACCATGCCTTT  
|| ||||| ||| | |||| | | | | | | |||  
CCGGACGAGAAGTCCAT---CACCTACGTGGTCACCTACTATCACTACTTT
```

Vlož pomlčky (medzery, gaps) tak, aby rovnaké bázy boli pod sebou.
Dobré zarovnanie má veľa zarovnaných rovnakých báz, málo medzier.

Na čo sú dobré zarovnanie?

- **Orientácia v obrovských databázach.**

Genbank má vyše 100 GB sekvencií.

Napr. odkiaľ z genómu je daná mRNA?

- **Určovanie funkcie (napr. proteínu).**

Podobné sekvencie často majú rovnakú/podobnú funkciu.

- **Štúdium evolúcie.**

Hľadáme homológy, sekvencie, ktoré sa vyvinuli z toho istého spoločného predka.

V ideálnom prípade medzery zodpovedajú inzerciam a deléciám, zarovnané bázy zachovaným bázam a substitúciám.

- **Hľadanie génov a iných funkčných prvkov.**

Menia sa pomalšie ako ostatné sekvencie.

Formulácia problému

Skórovanie zarovnaní: napr. zhoda +1, nezhoda -1, medzera -1.

```
GAGAAGGCCATAATGACCTATGTGTCCAGCT
|||||  |||  ||||  ||  ||  ||
GAGAAGTCCAT---CACCTACGTGGTCACCT
```

22 zhôd, 6 nezhôd, 3 medzery → skóre 13.

V praxi zložitejšie skórovanie. Chceme nastaviť tak, aby homológy mali vysoké skóre, náhodné zarovnanie nízke.

Problém 1: globálne zarovnanie (global alignment)

Vstup: sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$.

Výstup: zarovnanie X a Y s najvyšším skóre.

Problém 2: lokálne zarovnanie (local alignment)

Vstup: sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$.

Výstup: zarovnanie podreťazcov $x_i \dots x_j$ a $y_k \dots y_l$ s najvyšším skóre.

Dynamické programovanie pre globálne zarovnanie (Needleman, Wunsch 1970)

Podproblém: $A[i, j]$: najvyššie skóre globálneho zarovnaní reťazcov $x_1x_2 \dots x_i$ a $y_1y_2 \dots y_j$.

Jeden z reťazcov dĺžky 0: druhý reťazec je zarovnaný s medzerou.

$$A[0, j] = -j, \quad A[i, 0] = -i.$$

Všeobecný prípad, $i > 0, j > 0$:

ak $x_i = y_j$ a sú zarovnané $A[i, j] = A[i - 1, j - 1] + 1$,

ak $x_i \neq y_j$ a sú zarovnané $A[i, j] = A[i - 1, j - 1] - 1$,

ak x_i je zarovnané s medzerou $A[i, j] = A[i - 1, j] - 1$,

ak y_j je zarovnané s medzerou $A[i, j] = A[i, j - 1] - 1$.

Dynamické programovanie pre globálne zarovnanie

Podproblém: $A[i, j]$: najvyššie skóre globálneho zarovnanania reťazcov $x_1x_2 \dots x_i$ a $y_1y_2 \dots y_j$.

Všeobecný prípad, $i > 0, j > 0$:

ak x_i a y_j sú zarovnané $A[i, j] = A[i - 1, j - 1] + s(x_i, y_j)$

ak x_i je zarovnané s medzerou $A[i, j] = A[i - 1, j] - 1$

ak y_j je zarovnané s medzerou $A[i, j] = A[i, j - 1] - 1$

kde $s(x, y) = 1$ ak $x = y$ a $s(x, y) = -1$ ak $x \neq y$

Rekurencia:

$$A[i, j] = \max \begin{cases} A[i - 1, j - 1] + s(x_i, y_j), \\ A[i - 1, j] - 1, \\ A[i, j - 1] - 1 \end{cases}$$

Príklad globálneho zarovnaní

CATGTCATA vs CAGTCCTAGA

		C	A	G	T	C	C	T	A	G	A
	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
C	-1	1	0	-1	-2	-3	-4	-5	-6	-7	-8
A	-2	0	2	1	0	-1	-2	-3	-4	-5	-6
T	-3	-1	1	1	?						
G	-4										
T	-5										
C	-6										
G	-7										
T	-8										
A	-9										

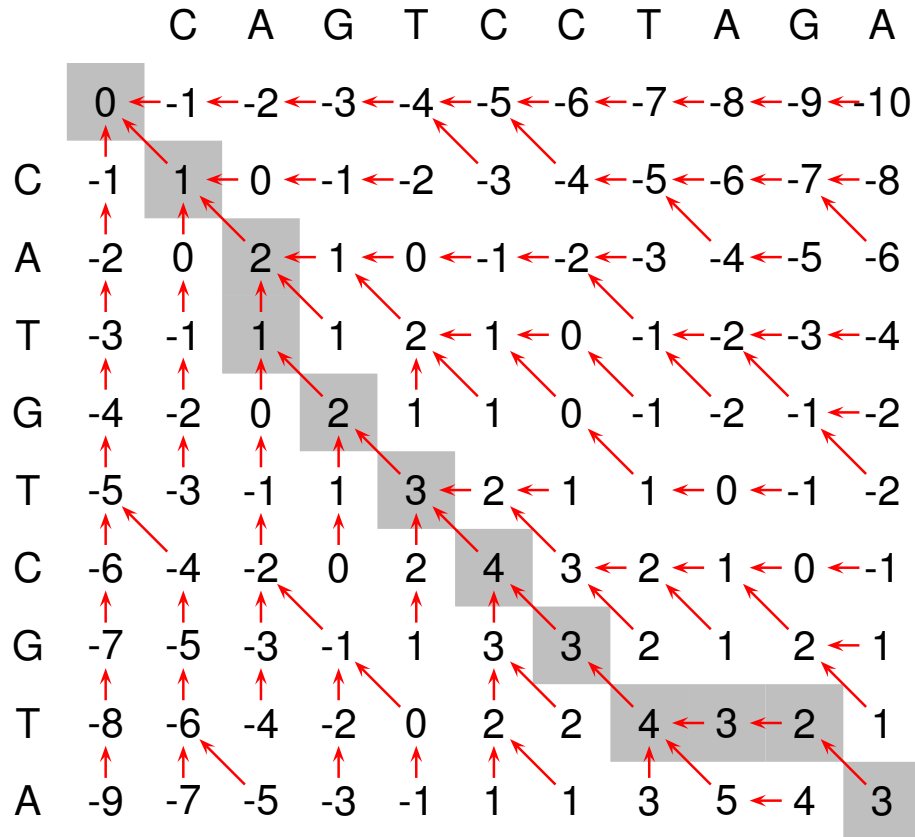
$$A[i, j] = \max \begin{cases} A[i - 1, j - 1] + s(x_i, y_j), \\ A[i - 1, j] - 1, \\ A[i, j - 1] - 1 \end{cases}$$

Príklad globálneho zarovnania

CATGTCATA vs CAGTCCTAGA

		C	A	G	T	C	C	T	A	G	A
	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
C	-1	1	0	-1	-2	-3	-4	-5	-6	-7	-8
A	-2	0	2	1	0	-1	-2	-3	-4	-5	-6
T	-3	-1	1	1	2	1	0	-1	-2	-3	-4
G	-4	-2	0	2	1	1	0	-1	-2	-1	-2
T	-5	-3	-1	1	3	2	1	1	0	-1	-2
C	-6	-4	-2	0	2	4	3	2	1	0	-1
G	-7	-5	-3	-1	1	3	3	2	1	2	1
T	-8	-6	-4	-2	0	2	2	4	3	2	1
A	-9	-7	-5	-3	-1	1	1	3	5	4	3

Ako získať zarovnanie?



CA-GTCCTAGA

CATGTCAT--A

Dynamické programovanie pre lokálne zarovnanie (Smith, Waterman 1981)

Podproblém: $A[i, j]$: najvyššie skóre lokálneho zarovnaní reťazcov $x_1x_2 \dots x_i$ a $y_1y_2 \dots y_j$, ktoré obsahuje bázy x_i a y_j , alebo je prázdne.

Jeden z reťazcov dĺžky 0: prázdne zarovnanie $A[0, j] = A[i, 0] = 0$

Všeobecný prípad, $i > 0, j > 0$:

ak x_i a y_j sú zarovnané $A[i, j] = A[i - 1, j - 1] + s(x_i, y_j)$

ak x_i je zarovnané s medzerou $A[i, j] = A[i - 1, j] - 1$

ak y_j je zarovnané s medzerou $A[i, j] = A[i, j - 1] - 1$

ak x_i a y_j nie sú časťou zarovnaní s kladným skóre $A[i, j] = 0$

Rekurencia: $A[i, j] = \max \left\{ \begin{array}{l} 0, \\ A[i - 1, j - 1] + s(x_i, y_j), \\ A[i - 1, j] - 1, \\ A[i, j - 1] - 1 \end{array} \right.$

Príklad lokálneho zarovnania

		C	A	G	T	C	C	T	A	G	A
	0	0	0	0	0	0	0	0	0	0	0
C	0	1	0	0	0	1	1	0	0	0	0
A	0	0	2	1	0	0	0	0	1	0	1
T	0	0	1	1	2	1	0	1	0	0	0
G	0	0	0	2	1	1	0	0	0	1	0
T	0	0	0	1	3	2	1	1	0	0	0
C	0	1	0	0	2	4	3	2	1	0	0
G	0	0	0	1	1	3	3	2	1	2	1
T	0	0	0	0	2	2	2	4	3	2	1
A	0	0	1	0	1	1	1	3	5	4	3

CA-GTCCTA

CATGTCATA

Zložitejšie skórovanie

Problémy +1, -1 skórovania:

- Je skutočne jedna nezhoda alebo medzera až taká zlá v porovnaní s jednou zhodou?
- Čo urobíme pre zarovnávanie proteínov?
(20 prvková abeceda \approx 200 parametrov)

Úloha skórovacej schémy:

- Chceme vedieť rozlíšiť **lepšie zarovnanie** od **horších zarovnaní**:
 - Ktoré usporiadania pomlčiek dávajú väčší zmysel
- Chceme vedieť, či dané zarovnanie **má biologický význam**:
 - Ide o homológy, alebo sekvencie nesúvisia?

Zložitejšie skórovanie: prvý pokus

Nech X a Y sú **správne zarovnané homológy**

a = pravdepodobnosť, že sa dve bázy **zhodujú**

b = pravdepodobnosť, že sa **nezhodujú**

c = pravdepodobnosť, že báza je **zarovnaná s medzerou**

$$a + b + c = 1$$

Pravdepodobnosť zarovnania A :

```
GAGAAGGCCATAATGACCTATGTGTCCAGCT
||||| |||  ||||| ||  ||  |
GAGAAGTCCAT---CACCTACGTGGTCACCT
```

$$\Pr(A) = a^{22}b^6c^3$$

Ktoré je pravdepodobnejšie?

```
CACA
|  |
CCAA
```

$$\Pr(A) = a^2b^2$$

```
CACA-
|  ||
C-CAA
```

$$\Pr(A) = a^3c^2$$

Zložitejšie skórovanie: prvý pokus

Zlogaritmujeme: násobenie sa zmení na sčítavanie
môžeme použiť S.-W. alebo N.-W. dyn. prog. algoritmy

$$\Pr(A) = a^{22}b^6c^3$$

$$\log \Pr(A) = 22 \log a + 6 \log b + 3 \log c$$

Skóre: Zhoda: $\log a$ Nezhoda: $\log b$ Medzera: $\log c$

Nevýhody takejto schémy:

- Vždy záporné skóre \Rightarrow čo s lokálnymi zarovnaniami?
- Neužitočné pre porovnávanie rôznych párov sekvencií

Zložitejšie skórovanie: dva pravdepodobnostné modely

(Pre jednoduchosť teraz neuvažujme medzery)

Model H: Sekvencie X a Y sú **správne zarovnané homológy**

$$\Pr(X, Y | H) = \prod_{i=1}^n p(x_i, y_i)$$

$p(x_i, y_i)$: pravdepodobnosť, že vidíme zarovnané práve bázy x_i a y_i

Model R: Sekvencie X a Y nijako spolu nesúvisia

$$\Pr(X, Y | R) = \left(\prod_{i=1}^n p(x_i)\right) \left(\prod_{i=1}^n p(y_i)\right)$$

$p(x_i)$: pravdepodobnosť výskytu bázy x_i

Porovnanie modelov H a R: “log likelihood”

$$\log \frac{\Pr(X, Y | H)}{\Pr(X, Y | R)}$$

Zložitejšie skórovanie: dva pravdepodobnostné modely

Porovnanie modelov H a R : “log likelihood”

$$\log \frac{\Pr(X, Y | H)}{\Pr(X, Y | R)}$$

- Dve sekvencie sú **homológy**
 - ⇒ pomer pravdepodobností je oveľa väčší ako 1
 - ⇒ **veľmi kladné skóre**
- Dve sekvencie **nesúvisia**
 - ⇒ pomer pravdepodobností je oveľa menší ako 1
 - ⇒ **veľmi zaporné skóre**

$$\log \frac{\Pr(X, Y | H)}{\Pr(X, Y | R)} = \log \frac{\prod_{i=1}^n p(x_i, y_i)}{(\prod_{i=1}^n p(x_i)) (\prod_{i=1}^n p(y_i))} = \sum_{i=1}^n \log \frac{p(x_i, y_i)}{p(x_i)p(y_i)}$$

BLOSUM62 skórovacia matica pre proteíny

BLOcks of aminoacid **S**Ubstitution **M**atrix; Henikoff, Henikoff 1992

	A	R	N	D	C	Q	E	G	H	I	L	...
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	
N	-2	0	6	1	-3	0	0	0	1	-3	-3	
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	
...												

- Vyber **biologicky relevantné zarovnanie** proteínov (BLOCKS)
- Páry s nanajvyš 62% identitou
- $p(x, y)$: ako často vidíme aminokyseliny x a y zarovnané
- $p(x)$: ako často sa vyskytuje aminokyselina x

- **skóre pre dvojicu aminokyselín x a y** : $\log \frac{p(x, y)}{p(x)p(y)}$
- prenásobíme konštantou a zaokrúhlime:
 - aby sme neurobili príliš veľkú chybu
 - aby sa s číslami lepšie počítalo

Zložitejšie skórovanie: afínne skóre medzier

```
CCCGACGAGAAGGCCATAATGACCTATGTGTCCAGCTTCTACCATGCCTTT
|| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
CCGGACGAGAAGTCCAT---CACCTACGTGGTCACCTACTATCACTACTTT
```

Niekoľko medzier za sebou asi nevzniklo nezávisle, možno jedna mutácia.

Penalta za začatie medzery (gap opening cost) o ,

Penalta za rozšírenie medzery o jedna (gap extension cost) e .

Medzera dĺžky g má penaltu $o + e(g - 1)$.

Zvolíme $o < e$ (t.j. $|o| > |e|$).

Základné nastavenia blastn: zhoda +2, nezhoda -3, $o = -5$, $e = -2$.

Príklad vyššie: 22 zhôd, 6 nezhôd, 1 medzera dĺžky 3

→ skóre $2 \cdot 22 - 3 \cdot 6 - 5 - 2 \cdot 2 = 16$.

Zhrnutie

- Globálne a lokálne zarovania
- Needleman-Wunschov a Smith-Watermanov algoritmus
- Skórovanie zarovnaní pomocou porovnávaní modelov
- Proteínové BLOSUM matice
- Afínne skórovanie medzier

Problémy na zamyslenie

1. **Časová zložitosť Smith-Waterman:** $O(nm)$

n - veľkosť prvej sekvencie

m - veľkosť druhej sekvencie

Čo robiť ak chceme porovnať ľudský genóm s myšacím genómom?

2. Povedzme, že nájdeme zarovnanie so skóre 14

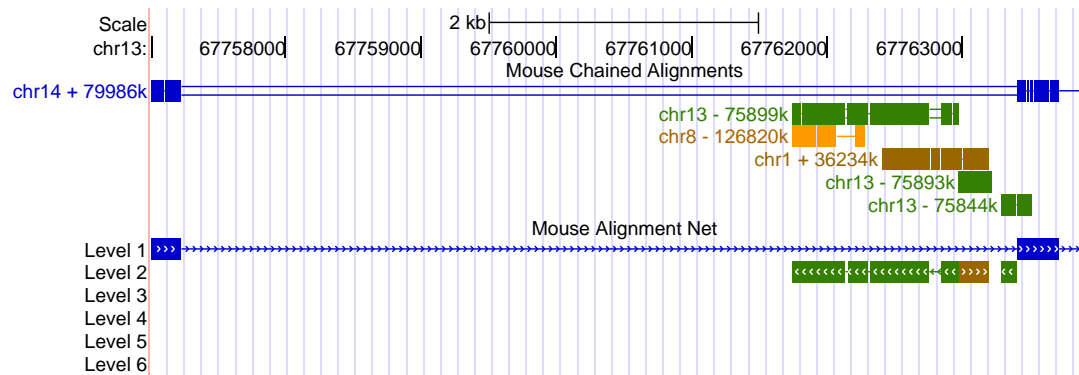
Je toto skóre dobré, alebo ide o niečo, čo vidíme náhodou?

Oznamy

- Výber článku na journal club formulárom na stránke do budúcej stredy 19.10. Zloženie skupín oznámime na budúcej prednáške.
- Domáca úloha 1 je zverejnená na stránke, odovzdávajújte do štvrtka 3.11. do 9:00 pod dvere M-163.
- Na domácich úlohách neopisujte. Môžete sa rozprávať, ale nerobte si pritom poznámky, neukazujte si navzájom svoje riešenia. Každý by mal napísať riešenie samostatne.
- Otázky k zadaniam a všeobecnú diskusiu k predmetu píšete do Facebookovej skupiny (môžete spolužiakom aj odpovedať), otázky k vašim riešeniam posielajte vyučujúcim e-mailom.

Zarovňavanie sekvencií 2/2 (sequence alignment)

Tomáš Vinar̄
13.10.2016



Zhrnutie z minulej prednášky

- **Problém globálneho a lokálneho zarovnaní**

Vstup: sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$.

Výstup: zarovnanie X a Y s najvyšším skóre

resp. zarovnaní podreťazcov $x_i \dots x_j$ a $y_k \dots y_\ell$ s najvyšším skóre.

- **Správny algoritmus na riešenie**

dynamické programovanie

- **Realistické skórovacie schémy**

Máme správny algoritmus na zarovnávanie, čo viac nám chýba?

Časová zložitosť: $O(nm)$ na sekvenciách dĺžky n a m .

Koľko je to času v skutočnosti?

(jednoduchá implementácia, náhodné sekvencie dĺžky n ,
bežný počítač)

n	čas výpočtu
100	0.0008s
1,000	0.08s
10,000	8s
100,000	13 minút (*)
1,000,000	22 hodín (*)
10,000,000	3 mesiace (*)
100,000,000	25 rokov (*)

Potrebujeme efektívnejší algoritmus,

najmä ak chceme pracovať s celými genómami

Heuristické lokálne zarovnávanie

- Nie je zaručené, že nájdeme najlepšie zarovnanie, ale program pobeží rýchlejšie.
- Prehľadá iba “sľubné” časti dyn. prog. matice.

Napríklad: BLASTN [Altschul et al., 1990],
FASTA [Pearson and Lipman, 1988]

- Nájdí krátke zhodujúce sa úseky dĺžky w (**jadrá zarovnaní**).
- Rozšír každé jadro pozdĺž uhlopriečky na zarovnanie bez medzier.
- Spoj zarovnaní na neďalekých uhlopriečkach medzerami.
- Lokálne vylepši zarovnanie dynamickým programovaním (možno vynechať).

Heuristické lokálne zarovnávanie

Príklad: začíname z jadier dĺžky $w = 2$
(V praxi sa používa $w = 10$ a viac.)

		C	A	G	T	C	C	T	A	G	A
	0	0	0	0	0	0	0	0	0	0	0
C	0	1	0	0	0	1	1	0	0	0	0
A	0	0	2	1	0	0	0	0	1	0	0
T	0	0	0	1	2	1	0	1	0	0	0
G	0	0	0	0	1	0	0	0	0	1	0
T	0	0	0	0	2	2	1	1	0	0	0
C	0	1	0	0	0	4	3	0	0	0	0
A	0	0	2	1	0	3	3	2	1	0	1
T	0	0	1	1	2	2	2	4	3	2	1
A	0	0	1	0	1	1	1	3	5	4	3

1. nájdi zhodné úseky
2. rozšír bez medzier
3. spoj medzerami

Ako nájdeme zhodujúce sa úseky?

- Vybudujeme “slovník” úsekov dĺžky w z prvej sekvencie.
- Nájdeme každý úsek z druhej sekvencie v slovníku.

Príklad: CAGTCCTAGA vs CATGTCATA

Slovník:

AG 2, 8

CA 1

CC 5

CT 6

GA 9

GT 3

TA 7

TC 4

Hľadaj:

CA → 1

AT → -

TG → -

GT → 3

TC → 4

CA → 1

AT → -

TA → 7

Rýchlosť heuristického algoritmu

Algoritmus:

- Nájdi jadrá zarovnaní (krátke zhodujúce sa úseky dĺžky w).
- **Drahý krok:** Rozširovanie/spájanie jadier do väčších zarovnaní.

Náhodné zhody dĺžky w : nie sú častou zarovnaní s vysokým skóre. Vyfiltrujeme ich pri rozširovaní, ale spomaľujú program.

Koľko náhodných zhôd?

Dva nukleotidy sa zhodujú s pravdepodobnosťou $1/4$.

w zhôd za sebou s pravdepodobnosťou 4^{-w} .

Stredná hodnota počtu zhôd $nm4^{-w}$.

Zvýšenie w o 1 zníži počet zhôd cca 4 krát.

Senzitivita heuristického algoritmu

Algoritmus:

- Nájdi jadrá zarovnaní (krátke zhodujúce sa úseky dĺžky w).
- **Drahý krok:** Rozširovanie/spájanie jadier do väčších zarovnaní.

Nenájdené zarovnaní: vysoké skóre, ale **nemajú jadro dĺžky w**

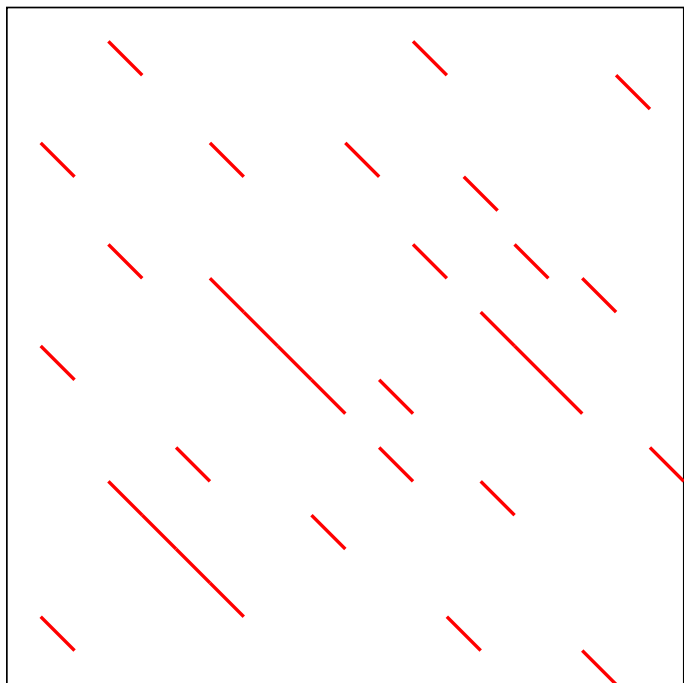
Príklad: CA-GTCCTA nenájdeme pre $w \geq 4$
 CATGTCATA

Senzitivita: aká časť **skutočných zarovnaní** obsahuje zhodu dĺžky w

Rýchlosť vs. senzitivita

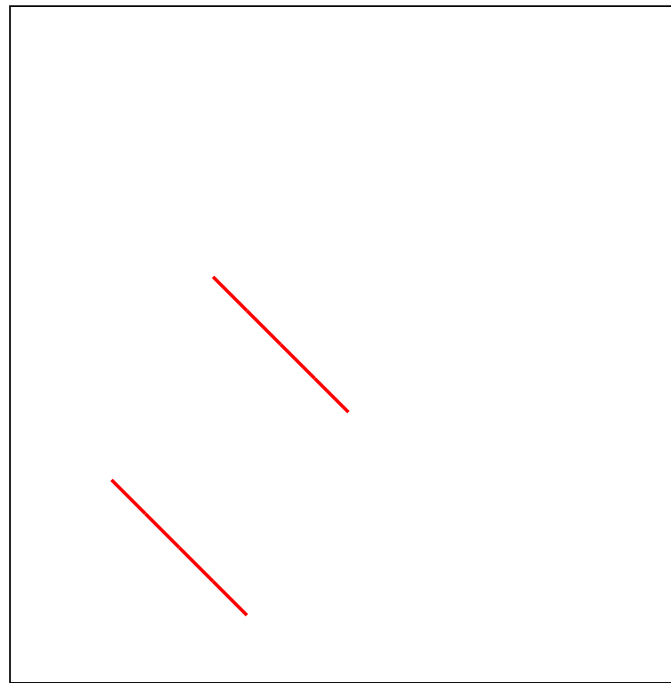
Malé w

veľa náhodných zhôd, pomalé



Veľké w

nenájdeme veľa zarovnaní



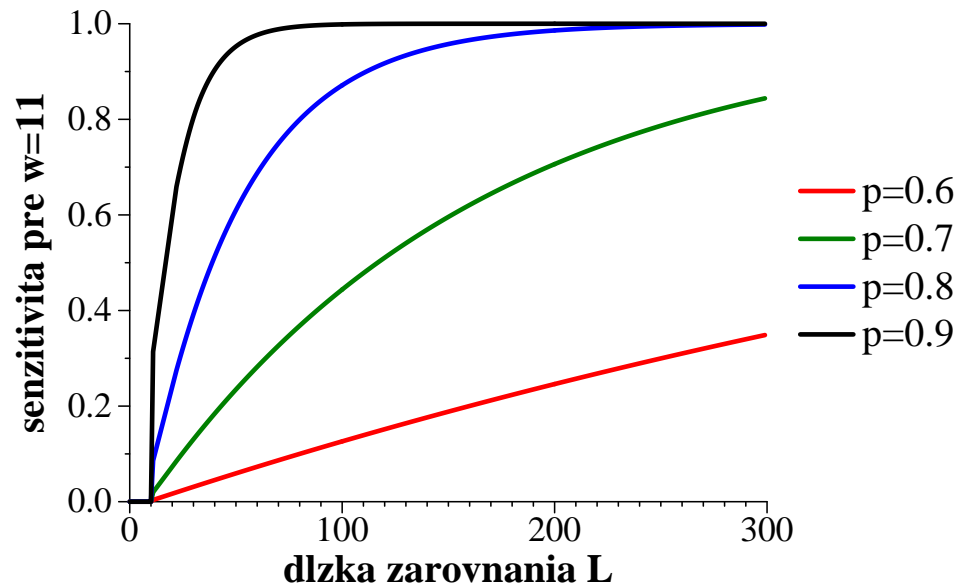
Senzitivita heuristického algoritmu

Odhad senzitivity:

Predpokladáme zarovnanie bez medzier, dĺžky L

Každá pozícia je zhoda s pravdepodobnosťou p

$f(L, p) = \Pr(\text{zarovnanie obsahuje } w \text{ zhôd za sebou})$



(človek-myš: $p \approx 0.7$)

BLAST algoritmus pre proteíny

BLOSUM62 skórovacia matica pre proteíny

	A	R	N	D	C	Q	E	G	H	I	...
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	
R	-1	5	0	-2	-3	1	0	-2	0	-3	
N	-2	0	6	1	-3	0	0	0	1	-3	
D	-2	-2	1	6	-3	0	2	-1	-1	-3	
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	
Q	-1	1	0	0	-3	5	2	-2	0	-3	
E	-1	0	0	2	-4	2	5	-2	0	-3	
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	
H	-2	0	1	-1	-3	0	0	-2	8	-3	
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	

Proteínový BLAST namiesto zhody dĺžky w vyžaduje 3 aminokyseliny so skóre aspoň 13

Áno: N I R
N L R
 $6+2+5=13$

Nie: A I L
A I L
 $4+4+4=12$

Príklady programov

NCBI BLAST: `blastn` pre DNA/RNA, `blastp` pre proteíny, `tblastx` preloží DNA do proteínu a použije `blastp`
[Altschul et al., 1990, Altschul et al., 1997]

UCSC Blat: veľmi rýchle vyhľadávanie veľmi podobných sekvencií, napr. mRNA ku genómu [Kent, 2002].

- používa veľké w
- vie rozdeliť mRNA na exóny

PSI-BLAST: [Altschul et al., 1997]

- Pre dotaz nájdeme zarovnanie cez `blastp`.
- Vidíme, ktoré pozície mutujú viac a ktoré menej.
- Nezhoda na zachovanej pozícii stojí viac.

⇒ nájde vzdialenejšie homológy.

Sequences producing significant alignments:		Score (Bits)	E Value	
ref XP_002345317.1 	PREDICTED: similar to protein tyrosine ph...	28.2	108	UG
ref XP_001726210.1 	PREDICTED: similar to protein tyrosine ph...	28.2	108	G
ref ZP_03264973.1 	isocitrate dehydrogenase, NADP-dependent [...	27.4	194	
ref XP_001225150.1 	hypothetical protein CHGG_07494 [Chaetomi...	27.4	194	G
ref YP_002967336.1 	hypothetical protein MexAM1_META2p1254 [M...	26.9	261	G
ref ZP_03013307.1 	hypothetical protein BACINT_00864 [Bactero...	26.9	261	
ref YP_001834672.1 	phospholipid/glycerol acyltransferase [Be...	26.9	261	G
ref ZP_04426281.1 	NADH dehydrogenase subunit L [Planctomyces...	26.1	469	
ref YP_003129642.1 	putative exonuclease RecJ [Halorhabdus ut...	26.1	469	G
ref ZP_02926313.1 	multidrug efflux pump, AcrB/AcrD/AcrF fami...	26.1	469	
ref ZP_02044690.1 	hypothetical protein ACTODO_01565 [Actinom...	26.1	469	
ref XP_001153320.1 	PREDICTED: similar to tyrosine phosphatas...	26.1	469	G
ref YP_001958968.1 	inner-membrane translocator [Chlorobium p...	26.1	469	GG
ref YP_003133865.1 	hypothetical protein Svir_20200 [Saccharo...	25.7	630	G

<http://blast.ncbi.nlm.nih.gov/Blast.cgi>

[Most Visited](#)
[Smart Bookmarks](#)
[Getting Started](#)
[Latest BBC Head...](#)
[Gmail](#)
[Entrez PubMed](#)

Alignments
[Select All](#)
[Get selected sequences](#)
[Distance tree of results](#)
[Multiple alignment](#) NEW

> [ref|XP_002345317.1|](#) **UG** PREDICTED: similar to protein tyrosine phosphatase 4a1 isoform 2 [Homo sapiens]
 Length=139

[GENE ID: 730167 LOC730167](#) | similar to protein tyrosine phosphatase 4a1 [Homo sapiens]

Score = 28.2 bits (59), Expect = 108
 Identities = 9/10 (90%), Positives = 10/10 (100%), Gaps = 0/10 (0%)

Query	1	VIVALASVEG	10
		V+VALASVEG	
Sbjct	79	VLVALASVEG	88

> [ref|XP_001726210.1|](#) **G** PREDICTED: similar to protein tyrosine phosphatase 4a1 isoform 1 [Homo sapiens]
 Length=170

[GENE ID: 730167 LOC730167](#) | similar to protein tyrosine phosphatase 4a1 [Homo sapiens]

Score = 28.2 bits (59), Expect = 108
 Identities = 9/10 (90%), Positives = 10/10 (100%), Gaps = 0/10 (0%)

Query	1	VIVALASVEG	10
		V+VALASVEG	
Sbjct	110	VLVALASVEG	119

Ako rozlíšiť, či ide o významné zarovnanie?

Zarovnanie so skóre S .

Dĺžka dotazu m . Veľkosť databázy n .

P -hodnota: Pravdepodobnosť, že pre náhodný dotaz dĺžky m v náhodnej databáze dĺžky n nájdeme zarovnanie so skóre aspoň S .

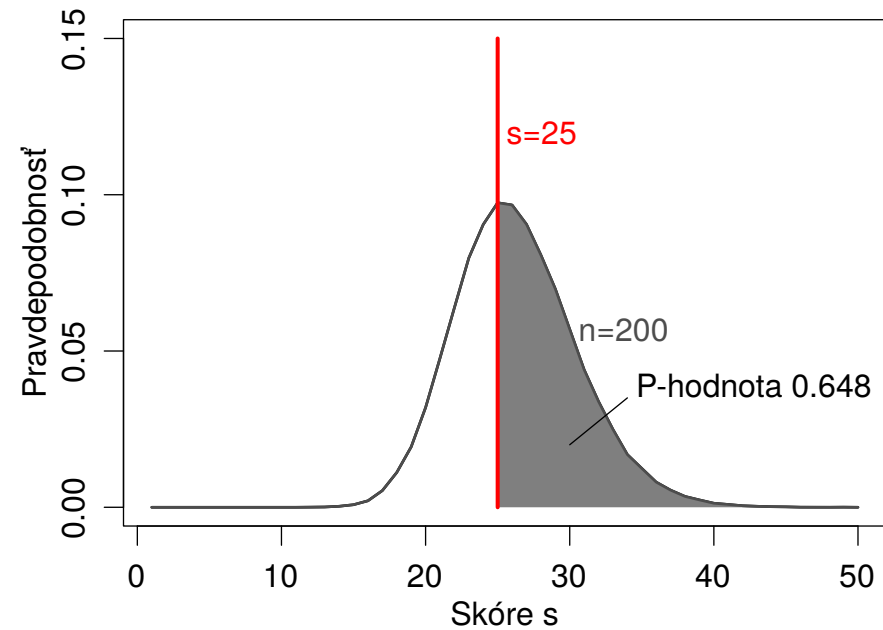
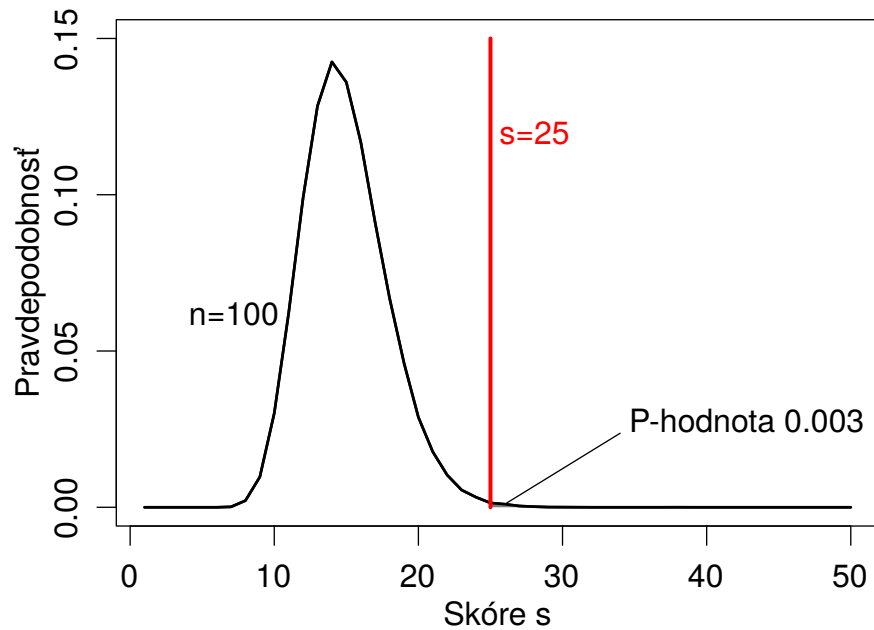
E -hodnota: Očakávaný počet zarovnaní so skóre aspoň S nájdených pre náhodný dotaz dĺžky m v náhodnej databáze dĺžky n .

Pri veľmi malých hodnotách sú E -value a P -value takmer identické.

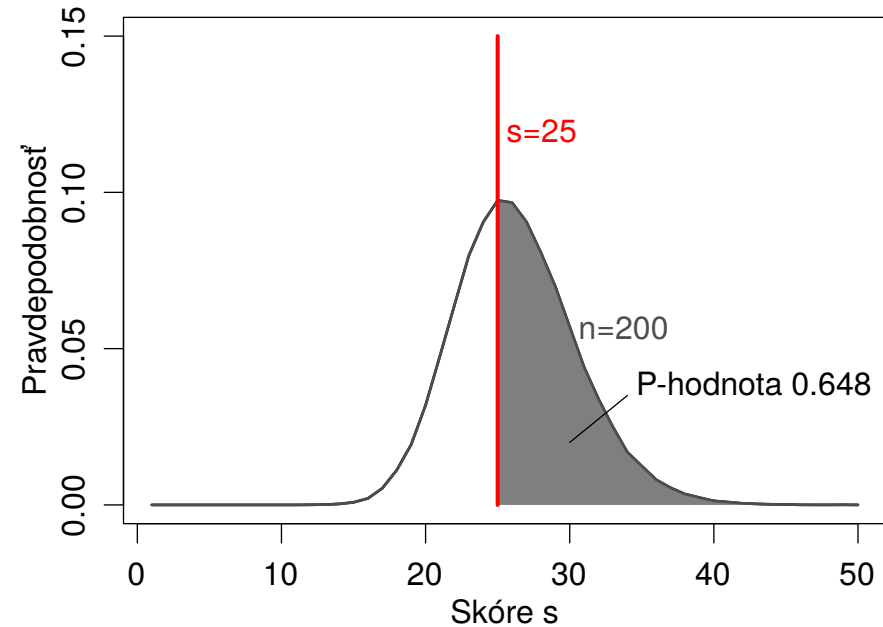
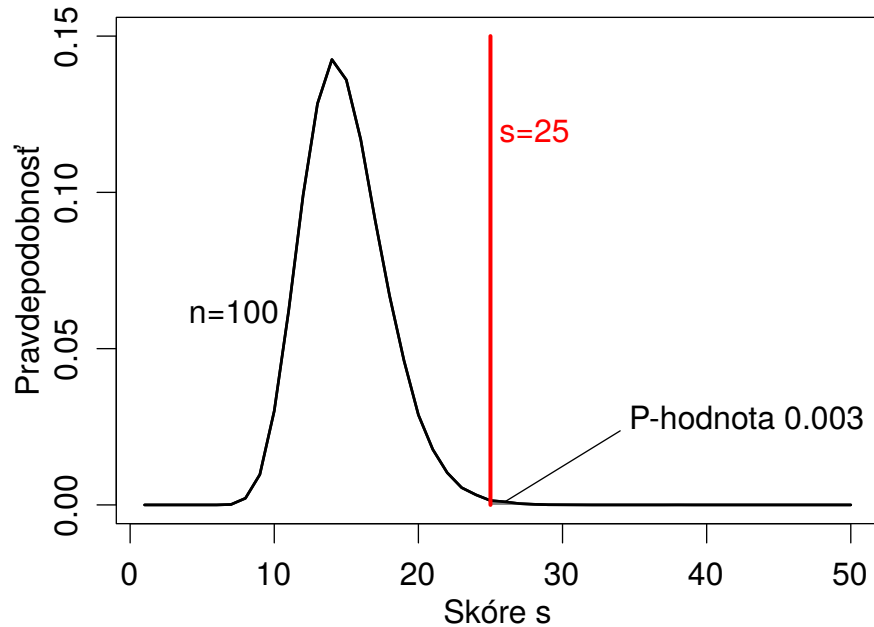
[Karlin and Altschul, 1990, Dembo et al., 1994]

Výpočet P-hodnoty simuláciou

- Vygenerujeme náhodne dve sekvencie dĺžky n
- Spočítame ich najlepšie lokálne zarovnanie (schéma +1/-1)
- Zaznamenáme si výsledné skóre
- Opakujeme veľa krát



Výpočet P-hodnoty simuláciou (pokr.)



P-hodnota pre skóre 25:

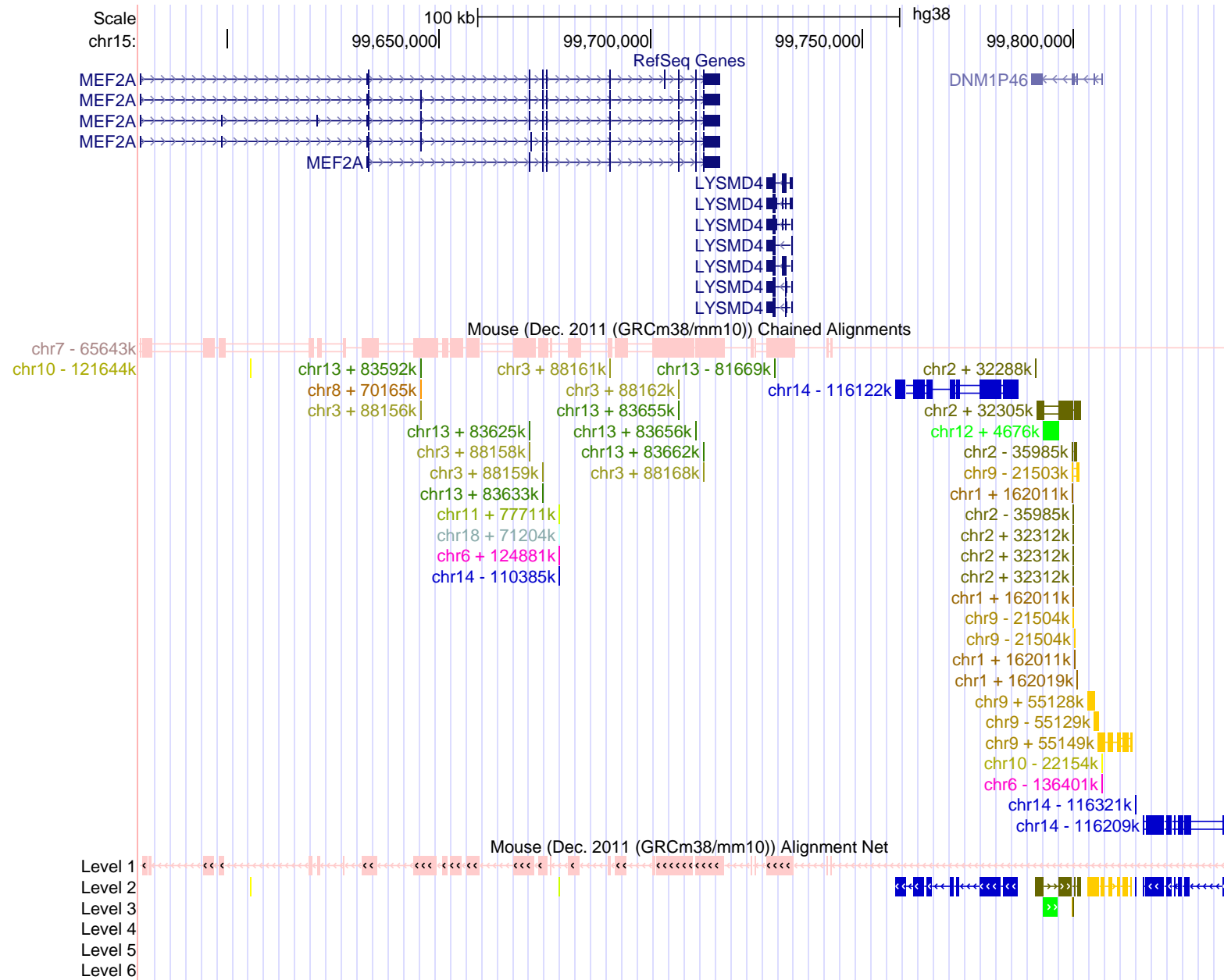
Aká časť zarovnaní má skóre 25 alebo vyššie?

(V praxi je simulácia pomalá, existujú odhady rozdelenia.)

Genomické zarovnanie (whole-genome alignments)

Ku každému úseku ľudského genómu nájsť zodpovedajúcu časť z myši, psa, sliepky, atď. (predpočítané v UCSC browseri) [Kent et al., 2003]

- Lokálne zarovnanie nájdu exóny a iné zachované časti, sú však úseky, ktoré sa príliš zmenili.
- Pri duplikovaných úsekoch nevieme rozhodnúť, ktoré dvojice úsekov patria k sebe.
- **Synténia (synteny)**: lokálne zarovnanie, ktoré sa nachádzajú v dvoch genómoch v tom istom poradí a orientácii.
Pomáha nám určiť, ktoré dvojice úsekov vznikli z tej istej oblasti v spoločnom predkovi (ortológ)



Viacnásobné zarovnanie, multiple sequence alignment

Zarovnaj viacero sekvencií.

Zložitosť: $O(2^k n^k)$ pre k sekvencií dĺžky n

Pre všeobecné k NP-ťažké.

```
Human  ctccatagcaatgt-cagagatagggcagagcggat-----ggtggtgac
Rhesus ctccatggcaatgt-cagagatagggcagagcggat-----gctggtgac
Mouse  ttt--tgacaaca--tagagac-tgagatagaaaat-----atgctgac
Dog    -tccccgctaatgtacaaagatggggcag-gaaga--a----tgtgctgaa
Horse  -tccacggcaatac-tggagatggggcagagcaga--agat-ggtgatgaa
Armadillo ctgcatagaaatct-cagagatgggggaaagcaga-----agacattcat
Opossum atccatggaaacat-cagaagtgggagaaatagaaga----tggcaatga-
Platypus acccggggaagggg-aagaggaagggccggccg-----
```

Heuristické algoritmy, napr. CLUSTAL-W [Higgins et al., 1996], MUSCLE [Edgar, 2004] a TBA [Blanchette et al., 2004].

Zhrnutie

- Zarovnávanie (alignment) je základný nástroj bioinformatiky
- Formulácia problému: voľba skórovacej schémy
- Riešenie problému: presné ale pomalé algoritmy a rýchlejšie heuristiky, ktoré nie vždy nájdu všetko
- Špecializované programy na rôzne úlohy súvisiace so zarovnávaním

Organizačné poznámky

- DÚ1 je zverejnená, odovzdávanie do 3.11. ráno
- Dnes na konci prednášky zverejníme rozdelenie skupín na journal club

Hľadanie génov

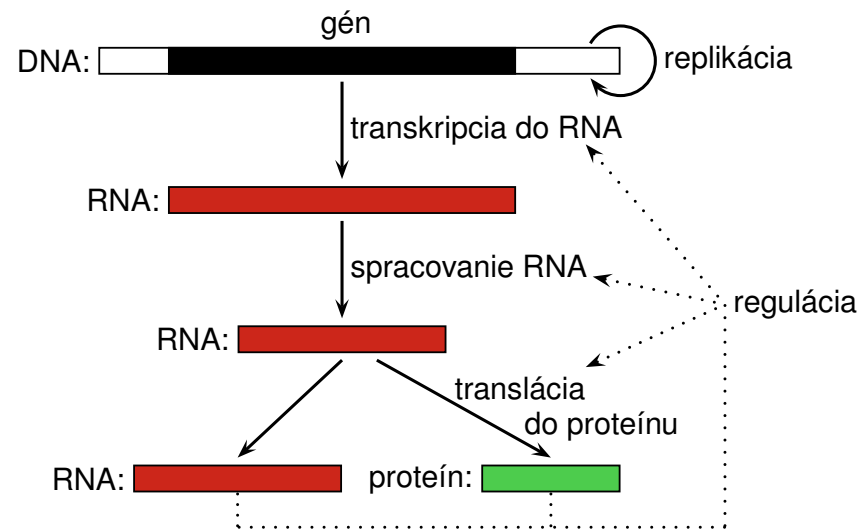
Broňa Brejová

20.10.2016

Čo s osekvenovanými genómami?

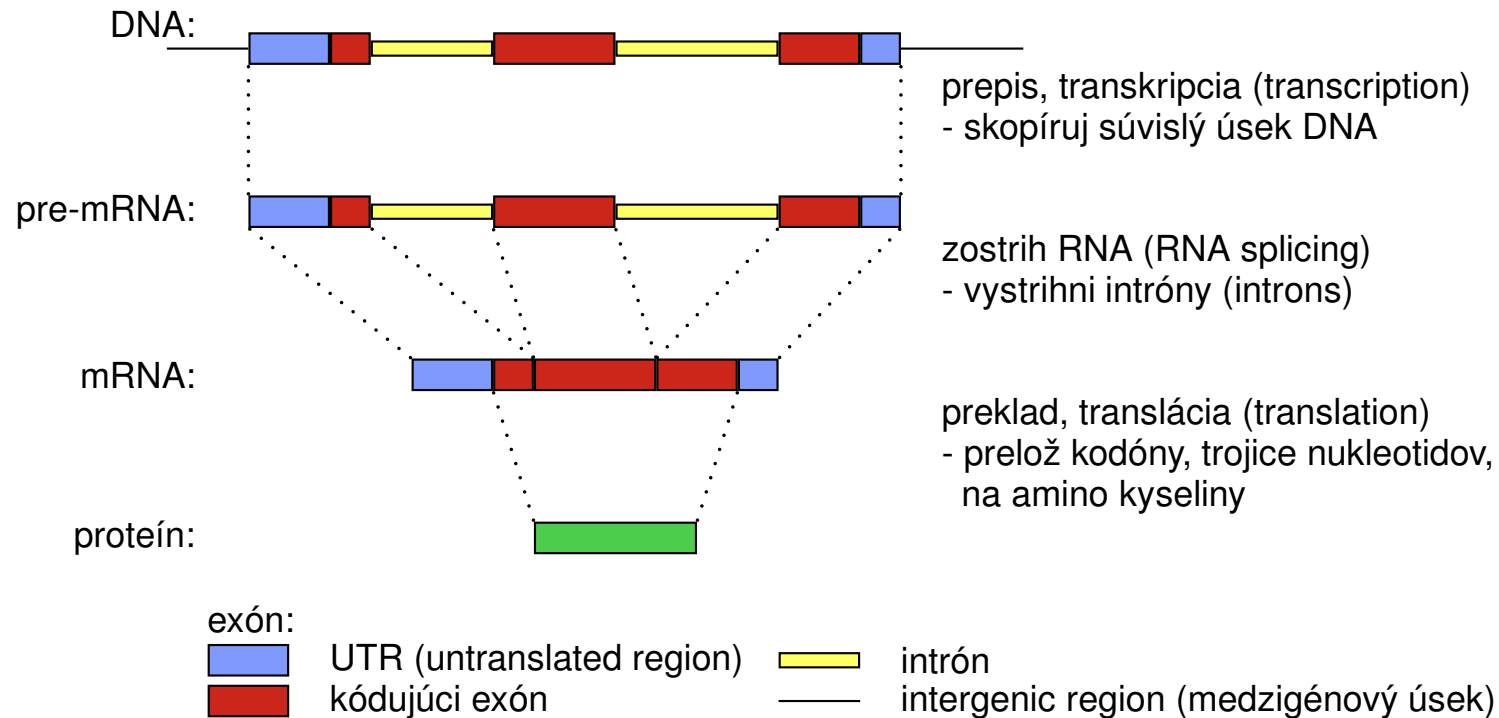
Chceme vedieť, čo genóm kóduje, hľadáme zaujímavé prvky, ako:

- gény kódujúce proteíny (dnešná prednáška)
- RNA gény
- signály pre reguláciu transkripcie, zostrihu, atď
- pseudogény (nefunkčné kópie génov)
- repetitívne sekvencie, opakovania (sequence repeats)

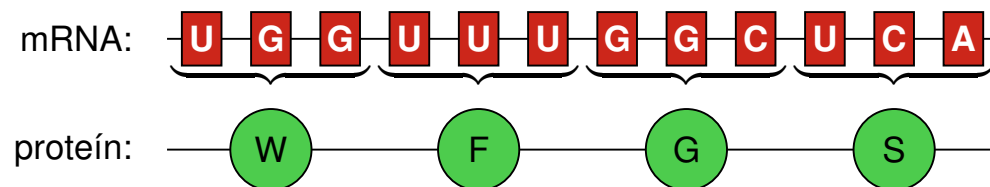


Štruktúra eukaryotických génov

Proces tvorby proteínov:



Translácia: tri bázy mRNA (kodón) → aminokyselina proteínu



Ľudský genóm

- gény kódujúce proteíny
 - cca 20,000, pokrývajú 40% genómu
 - cca 10 exónov v géne
 - exóny pokrývajú 2% genómu
 - kódujúce exóny 1.2% genómu

- repetitívne sekvencie
 - pokrývajú 49% genómu

Bioinformatický problém: hľadanie génov

Cieľ: nájsť všetky gény kódujúce proteíny v genóme.

Tým získame katalóg všetkých proteínov.

Zjednodušená:

- neuvažujeme alternatívny zostrih, prekrývajúce sa gény
- nehľadáme neprekladané oblasti (UTRs) na začiatku a konci génu

Bioinformatický problém: hľadanie génov

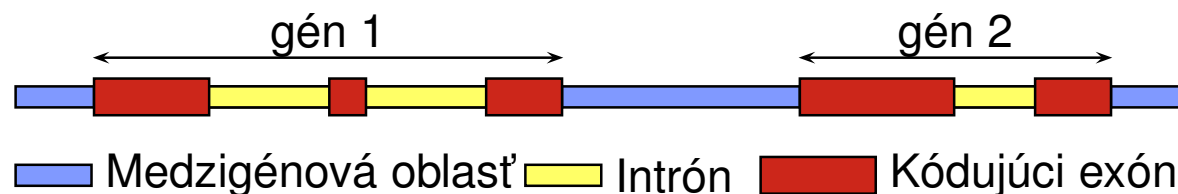
Vstup: sekvencia DNA

```
cggtgaaactgcacgattggtgctggcttaaagatagaccaatcagagtgtgtaacgtca  
tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca  
tgggcgtatattgcgctagtgttgggtgttccgctgtgctgtttttccgtcatggctcgca  
ctaagcaaactgctcgggaagtctactggtggcaaggcgccacgcaaacagttggccacta  
aggcagcccgcaaaagcgctccggccaccggcggcgtgaaaaagccccaccgctaccggc  
cgggcaccgtggctctgcgcgagatccgccgttatcagaagtccactgaactgcttattc  
gtaaactacctttccagcgcctggtgcgcgagattgcgcaggactttaaacagacctgc  
gtttccagagctccgctgtgatggctctgcaggaggcgtgcgaggcctacttggtagggc  
tatttgaggacactaacctgtgcgccatccacgccaagcgcgtcactatcatgccaagg  
acatccagctcgcccgccatccgcggagagagggcgtgattactgtggtctctctgac
```

Bioinformatický problém: Hľadanie génov

Cieľ: označ každú bázu ako intrón/exón/medzigénovú oblasť

```
cgggtgaaactgcacgattggtgctggcttaaagatagaccaatcagagtgtgtaacgtca
tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca
tgggcgtatTTGCGctagtgttgggtgttccgctgtgctgtttttccgtcatggctcgca
ctaagcaaactgctcggaaagtctactggtggcaaggcgccacgcaaacagttggccacta
aggcagcccgcaaaagcgcctccggccaccggcggcgtgaaaaagccccaccgctaccggc
cgggcaccgtggctctgcgcgagatccgccgttatcagaagtccactgaactgcttattc
gtaaactacctttccagcgcctggtgcgcgagattgcgcaggactttaaaacagacctgc
gtttccagagctccgctgtgatggctctgcaggaggcgtgcgaggcctacttggtagggc
tatttgaggacactaacctgtgcgccatccacgccaagcgcgtcactatcatgccaagg
acatccagctcgcccgcgcatccgcgagagagggcgtgattactgtggtctctctgac
```



Bioinformatický problém: hľadanie génov

Vstup: sekvencia DNA

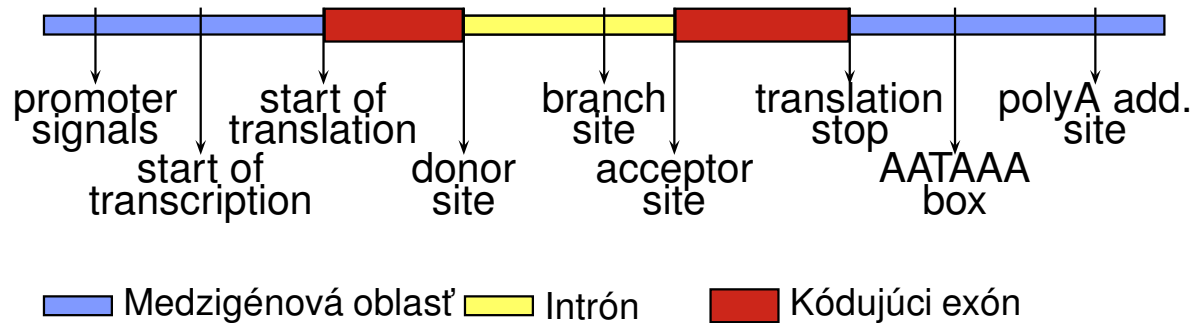
Cieľ: označ každú bázu ako intrón/exón/medzigénovú oblasť
(anotácia)

- Toto nie je dobre definovaný problém!
Ako spoznáme, čo je gén?

Ako spoznáme gény?

Signály na hraniciach exónov:

krátke reťazce, kde sa viažu komplexy zúčastňujúce sa na expresii génu



Príklad signálu: miesto zostrihu

Exón **Intrón**

```
ccatcccctatatatttatggcagGTgaggaaaggggtgggggctgggg
attcatcatcatgggtgcatcgGTgagtatctccaggccccaatc
agaagatctacccaccatctgGTAagtgtgtcccaccactgcccc
acagagtgagcccttcttcaagGTgggtgggtgtcagggcctcccc
acgagtcctgcatgagccagatGTAaggcttgccgttgccctcct
tgcagaacctcatgggtgctgagGTggggccaagcctgggcccggggg
tcgatgaatttgggatcatccgGTgagagctcttctctctctctgg
agatgacgtccgtgatgagaagGTaggggggtgcaccccagtccca
gtggagaatgagaggtgggatgGTaggtgatgccttcgaggccag
tttcttgtggcctattttaaaagGTAattcatggagaaatagaaaa
```

Ako spoznáme gény?

Zloženie sekvencie:

- iná frekvencia k -tic báz v kódujúcich a nekódujúcich oblastiach,
- kódujúce oblasti sú 3-periodické,
- stop kodóny (TAA, TGA, TAG) len na konci posledného kódujúceho exónu.

Príklad: ak uvažujeme len jednotlivé bázy, exóny majú viac C a G (ľudský genóm)

		a	c	g	t
kódujúci exón	0	0.26	0.26	0.32	0.16
	1	0.30	0.24	0.20	0.26
	2	0.17	0.32	0.31	0.20
intrón		0.26	0.22	0.22	0.30
medzig.		0.27	0.23	0.23	0.27

Bioinformatický problém: hľadanie génov

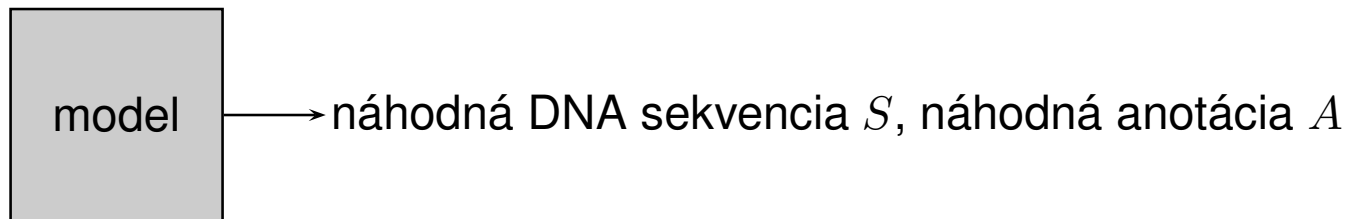
Vstup: sekvencia DNA

Cieľ: označ každú bázu ako intrón/exón/medzigénovú oblasť (anotácia)

- Toto nie je dobre definovaný problém!
Ako spoznáme, čo je gén?
- Žiadna informácia nám neumožňuje jednoznačne určiť, čo je gén.
- Chceme **skórovací systém**, ktorý povie, ako dobre potenciálna anotácia zodpovedá našim znalostiam.
- Potom hľadáme anotáciu (sadu neprekrývajúcich sa génov) **s maximálnym skóre.**
- Na definíciu skórovacieho systému použijeme **pravdepodobnostné modely.**

Pravdepodobnostný model génov

Žiadna informácia nám neumožňuje jednoznačne určiť, čo je gén.
Skombinujeme dostupnú informáciu pravdepodobnostným modelom.



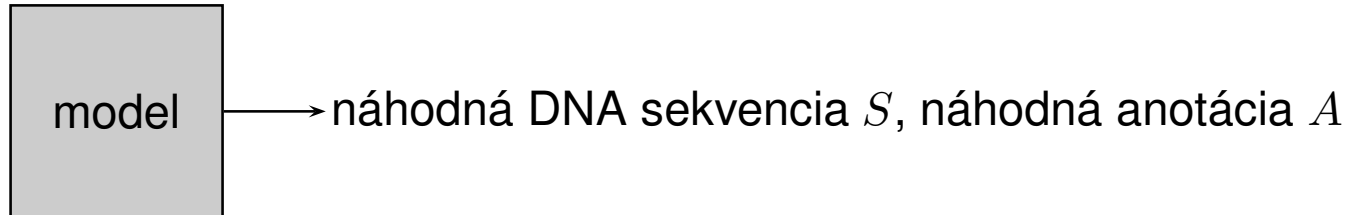
$\Pr(S, A)$ – pravdepodobnosť, že model vygeneruje pár (S, A) .

Model zostavíme tak, aby páry s vlastnosťami podobnými skutočným génom mali veľkú pravdepodobnosť.

Použitie: pre novú sekvenciu S nájdí najpravdepodobnejšiu anotáciu

$$A = \arg \max_A \Pr(A|S)$$

Pravdepodobnostný model génov



Použitie: pre sekvenciu S nájdí najpravdepodobnejšiu anotáciu A

Hračkársky príklad modelu: sekvencie dĺžky 2

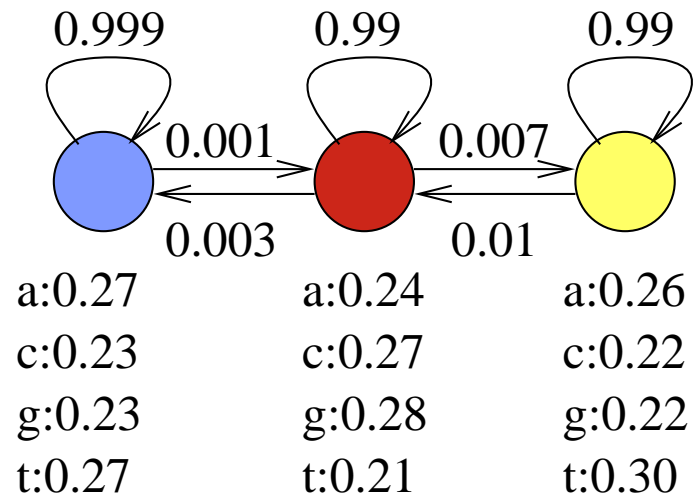
Tabuľka pravdepodobností pre 16 sekvencií, 9 anotácií (súčet 1)

Najpravdepodobnejšia anotácia pre $S = aa$ je **aa**.

aa	0.008	ac	0.009	ag	0.0085	...
aa	0	ac	0	...		
aa	0.011	...				
aa	0					
aa	0.009					
aa	0					
aa	0.007					
aa	0					
aa	0.010					

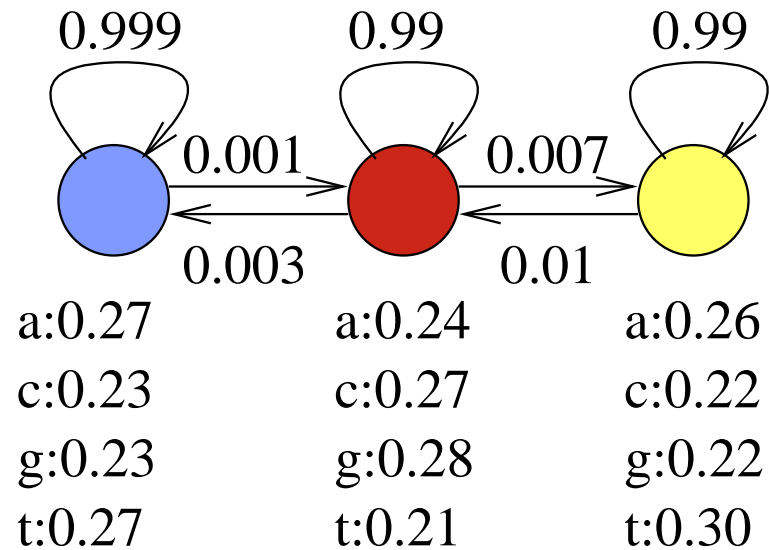
Skrytý Markovov model, hidden Markov model (HMM)

Spôsob, ako zdefinovať model pre dlhšie sekvencie.



- Konečný automat, stavy napr. exón, intrón, medzigénová oblasť
- Sekvenciu aj anotáciu generuje bázu po báze
- V každom kroku je v jednom stave a náhodne vygeneruje jednu bázu podľa tabuľky v stave
- Potom sa presunie do ďalšieho stavu podľa pravdepodobností na hranách

Skrytý Markovov model (HMM)



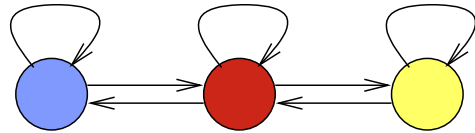
Předpokladajme, že model vždy začíná v modrom stave.

Príklad:

$$\Pr(\text{a} \color{red}{\text{c}} \text{a}) = 0.27 \cdot 0.001 \cdot 0.27 \cdot 0.99 \cdot 0.24 = 0.000017$$

$$\Pr(\text{a} \color{blue}{\text{a}} \text{a}) = 0.27 \cdot 0.999 \cdot 0.23 \cdot 0.999 \cdot 0.27 = 0.017$$

Matematické označenie



Sekvencia S_1, \dots, S_n


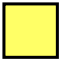



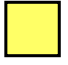
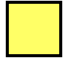

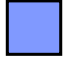
Anotácia A_1, \dots, A_n

Parametre modelu:

Prechodová pravdepodobnosť $a(u, v) = \Pr(A_{i+1} = v | A_i = u)$,

Emisná pravdepodobnosť $e(u, x) = \Pr(S_i = x | A_i = u)$,

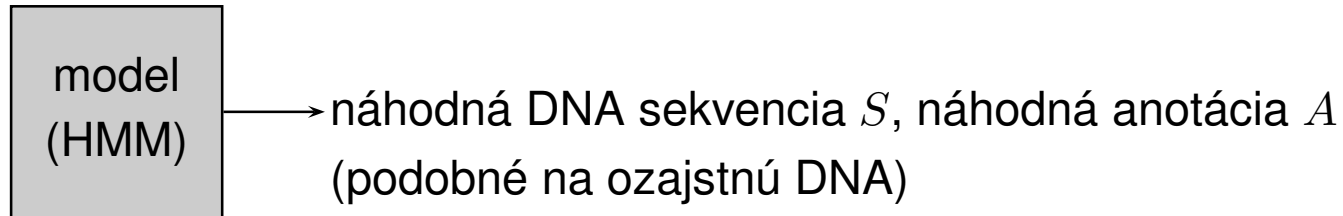
Počiatočná pravdepodobnosť $\pi(u) = \Pr(A_1 = u)$.

a				e	a	c	g	t
	0.99	0.007	0.003		0.24	0.27	0.28	0.21
	0.01	0.99	0		0.26	0.22	0.22	0.30
	0.001	0	0.999		0.27	0.23	0.23	0.27

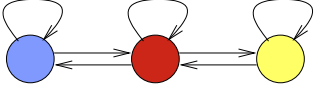
Výsledná pravdepodobnosť: $\Pr(A_1, \dots, A_n, S_1, \dots, S_n) =$

$$\pi(A_1)e(A_1, S_1) \prod_{i=2}^n a(A_{i-1}, A_i)e(A_i, S_i)$$

Hľadanie génov s HMM



$\Pr(S, A)$ – pravdepodobnosť, že model vygeneruje pár (S, A) .

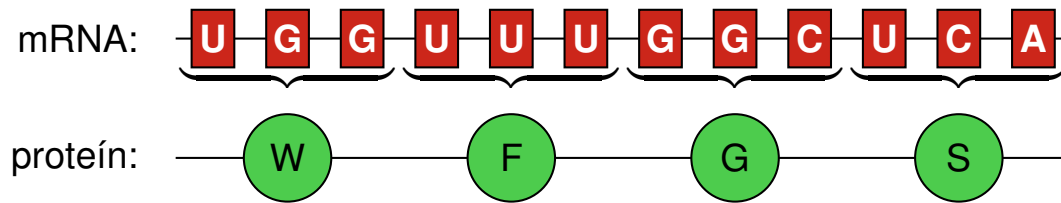
- **Určenie stavov a prechodov v modeli:** ručne, na základe poznatkov o štruktúre génu. 

- **Trénovanie parametrov:** emisné a prechodové pravdepodobnosti určíme na základe sekvencií so známymi génmi (**trénovacia množina**).

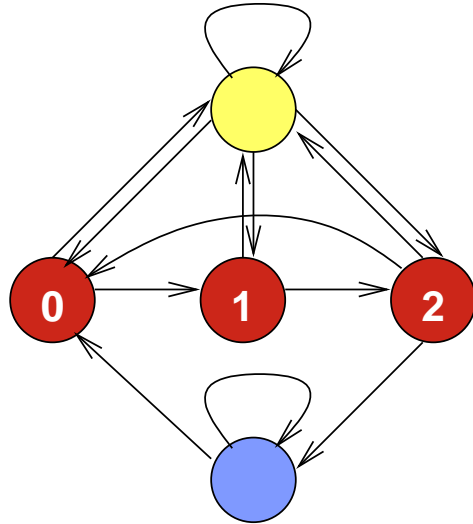
- **Použitie:** pre novú sekvenciu S nájdí najpravdepodobnejšiu anotáciu $A = \arg \max_A \Pr(A|S)$
Viterbiho algoritmus v čase $O(nm^2)$ (dynamické programovanie)

HMM na hľadanie génov: 3-periodické exóny

Kodón (trojica báz) → jedna aminokyselina



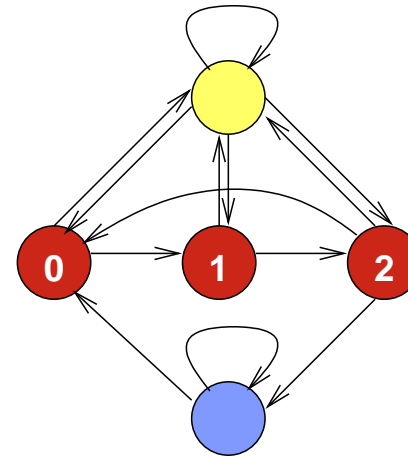
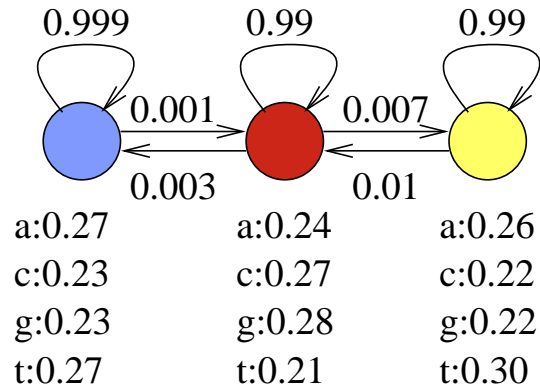
Namiesto jedného stavu pre exón použijeme tri stavy v cykle.



a	0	1	2	Yellow	Blue
0	0		0		0
1	0	0			0
2		0	0		
Yellow					0
Blue		0	0	0	

$\Pr(A_i|A_{i-1})$

Nové stavy mají odlišné emisné pravdepodobnosti

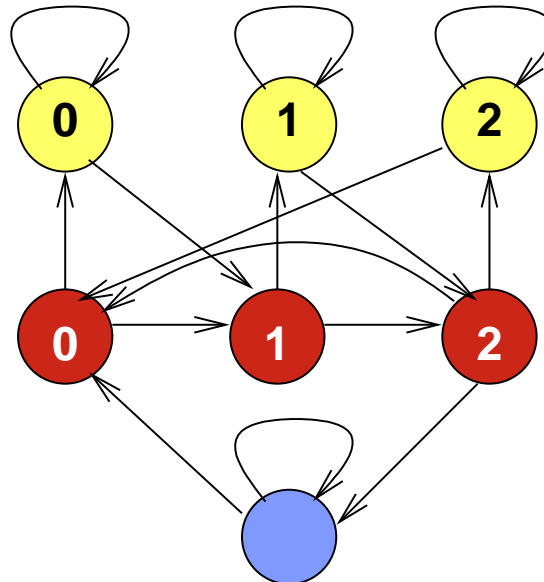
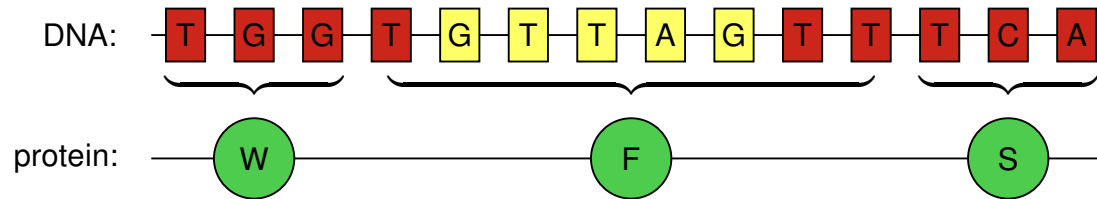


<i>e</i>	a	c	g	t
■	0.24	0.27	0.28	0.21
■	0.26	0.22	0.22	0.30
■	0.27	0.23	0.23	0.27

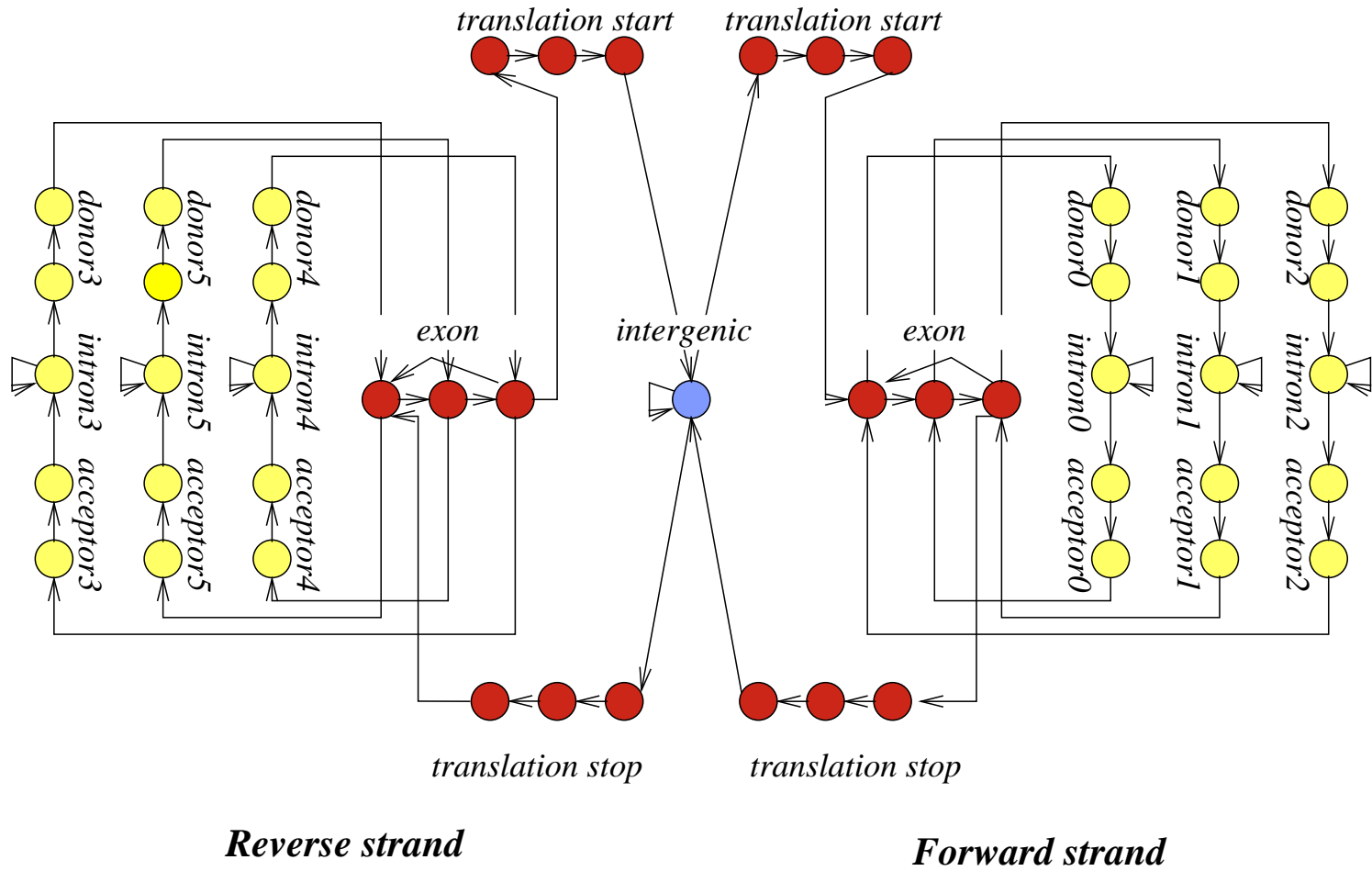
<i>e</i>	a	c	g	t
0	0.26	0.26	0.32	0.16
1	0.30	0.24	0.20	0.26
2	0.17	0.32	0.31	0.20
■	0.26	0.22	0.22	0.30
■	0.27	0.23	0.23	0.27

HMM na hľadanie génov: konzistentné kodóny

Intrón môže prerušiť kodón uprostred, chceme pokračovať, kde sme prestali.





HMM na hľadanie génov: celkový model



Stavy vyšších rádov

Rád 0: emisná tabuľka e určuje $\Pr(S_i|A_i)$

Rád 1: e určuje $\Pr(S_i|A_i, S_{i-1})$

A_i	S_{i-1}	a	c	g	t
	a	0.24	0.23	0.34	0.19
	c	0.30	0.31	0.13	0.26
	g	0.27	0.28	0.28	0.17
	t	0.13	0.28	0.38	0.21
	a	0.30	0.18	0.27	0.25
	c	0.32	0.28	0.06	0.35
	g	0.27	0.22	0.27	0.24
	t	0.20	0.21	0.26	0.33

...

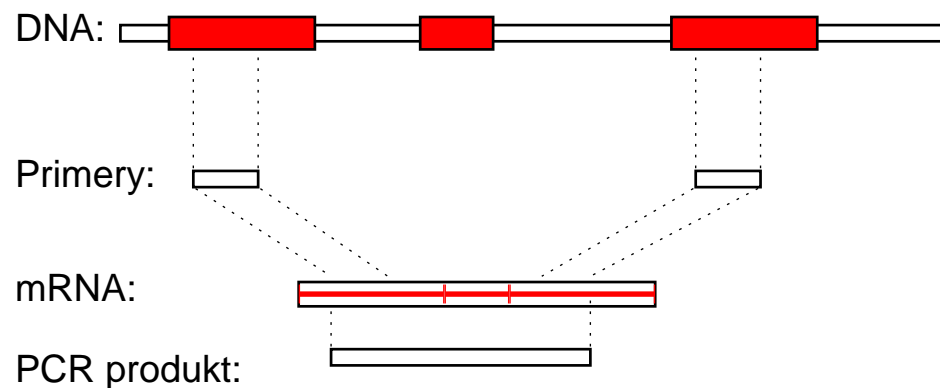
Na charakterizovanie exónov, intrónov atď používame rád 4-5.

Experimentálne overovanie génov

Overenie transkripcie a zstrihu

- RNA-Seq: sekvenovanie častí mRNA extrahovaných z bunky. Nie je cieleňé na konkrétny gén.
- RT PCR: cieleňe over konkrétny predpovedaný gén pomocou špecifických primerov.

Problémy: ťažko nájsť gény s expresiou iba za zvláštnych podmienok, napr. v embryu, kontaminácia genómovou DNA, nejednoznačné namapovanie na genóm.



Experimentálne overovanie génov

Overenie translácie, prítomnosti proteínu

- Hmotnostná spektrometria (mass spectrometry) dokáže detekovať prítomnosť proteínu izolovaného napr. z 2D gélu.
- Metódy založené na protilátkach (antibody), prípadne špecifické techniky podľa typu proteínu.

Príklady programov na hľadanie génov

Len na základe sekvencie DNA:

HMMGene [Krogh, 1997] (autor je priekopníkom HMM v bioinf.),
Genscan [Burge and Karlin, 1997] (po mnohé roky štandard),
GeneZilla [Majoros et al., 2004], ExonHunter [Brejová et al., 2005],
Augustus [Stanke and Waack, 2003] (novšie programy založené na
zovšeobecnených HMM).

CONTRAST [Gross et al., 2007], CONRAD [DeCaprio et al., 2007]
(najnovšia generácia založená na conditional random fields)

Prokaryotické genómy:

GeneMark [Lukashin and Borodovsky, 1998], Glimmer
[Delcher et al., 1999] a ďalšie.

Vybrané programy na hľadanie génov

Porovnávaním viacerých sekvencií:

Twinscan [Korf et al., 2001]

(prvý úspešný gene finder s dvoma genómami),

Exoniphy [Siepel and Haussler, 2004]

(viacero genómov, nehľadá celé gény),

N-SCAN [Gross and Brent, 2006]

(rozšírenie Twinscanu na viacero genómov).

Iná informácia: (napr. RNA-seq, príbuzné proteíny a pod.)

ExonHunter [Brejová et al., 2005], Augustus [Stanke et al., 2006],

Jigsaw [Allen and Salzberg, 2005],

Fgenesh++ [Solovyev et al., 2006].

Obmedzenia hľadačov génov

- Alternatívny zostrih (alternative splicing): jeden gén môže vyprodukovať viacero mRNA molekúl. Programy väčšinou hľadajú iba jednu.

Retained intron:



Skipped exon:



Alternative donor or acceptor:

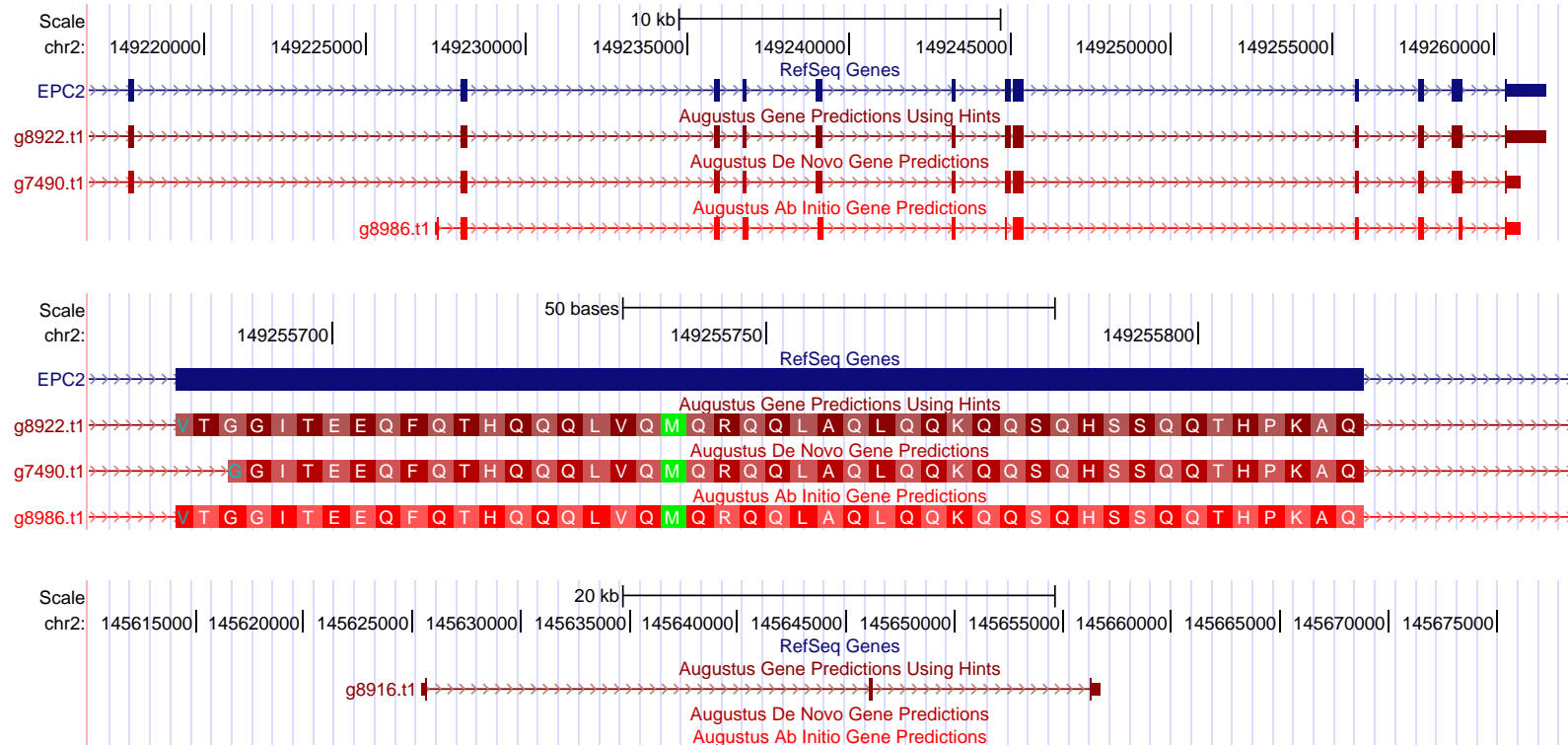


Mutually exclusive exons:



- Pretínajúce sa gény, resp. gény v intrónoch.
- Netypické gény (neobvyklé signály, veľmi krátke alebo dlhé exóny alebo intróny atď.)
- Hľadanie UTR a začiatku/konca transkripcie.

Hľadáče génov robia často chyby



Najlepšie metódy v 2005 na ľudskom genóme: [Guigo et al 2006]

20% génov, 60% exónov správne iba na základe DNA

35% génov, 65% exónov správne komparatívne

70% génov, 85% exónov správne s ďalšou informáciou

Koľko g3nov m3 3lovek?

Do 2001: R3zne odhady: **50 000–140 000** g3nov

2001: predbeŹn3 verzia ľudsk3ho gen3mu: **30 000–40 000** g3nov

2004: sekvencia ľudsk3ho gen3mu: **20 000–25 000** g3nov

2007: v katal3goch Ensembl, RefSeq a VEGA spolu **24 500** g3nov

[Clamp a kol. 2007] tvrdia, Źe iba **20 500** z nich je spr3vn3ch

Ale s3 g3ny, o ktor3ch eŹte nevieme?

2010: RefSeq m3 **22 333** g3nov

St3le neistota ± 1000 [Pertea, Salzberg 2010]

R3zni ľudia sa m3Źu l3iŹ v desiatkach g3nov

2012: Projekt ENCODE odhaduje **20 687** g3nov k3duj3cich prote3ny,

v priemere 6 altern3vn3ch transkriptov na g3n,

plus 8 800 kr3tk3ch a 9 600 dlh3ch RNA g3nov

Zhrnutie

- Novo osekvenované genómy treba anotovať:
určovať funkcie jednotlivým oblastiam sekvencie
- Príkladom anotácie je hľadanie génov kódujúcich proteíny
- Na hľadanie génov sa hodia skryté Markovove modely
- Modely robia veľa chýb, ale dajú nám základnú predstavu o polohe a počte génov, môžeme študovať ich funkciu

Journal club

- Vyhlásime rozdelenie do skupín, každá skupiny sa zoznámte, vymeňte si e-maily.
- Každý si najprv prečíta článok, potom sa koná stretnutie, kde o článku diskutujete, vysvetlíte si navzájom nejasnosti, plánujete písanie správy
- Prvé stretnutie skupiny najneskôr 17.11. (na FMFI alebo PriFUK), čas a miesto oznámte aspoň 2 dni vopred na facebookovej skupine predmetu
- Po stretnutí pošlite e-mail B.Brejovej s krátkou správou zo stretnutia
- Ak treba, dohodnite si s nami konzultácie

Správa zo journal clubu

- Vlastnými slovami hlavné metódy a výsledky článku
- Pochopiteľná pre študentov tohto predmetu (inf aj bio)
- Netreba pokryť všetko a naopak, môžete využiť aj iné zdroje
- Skúste vložiť vlastný pohľad na tému
- Rozsah cca 1-2 strany na osobu, jeden ucelený text
- V správe vymenujte členov skupiny, ktorí sa podieľali na jej spísaní, dostanú rovnako bodov

Organizačné poznámky

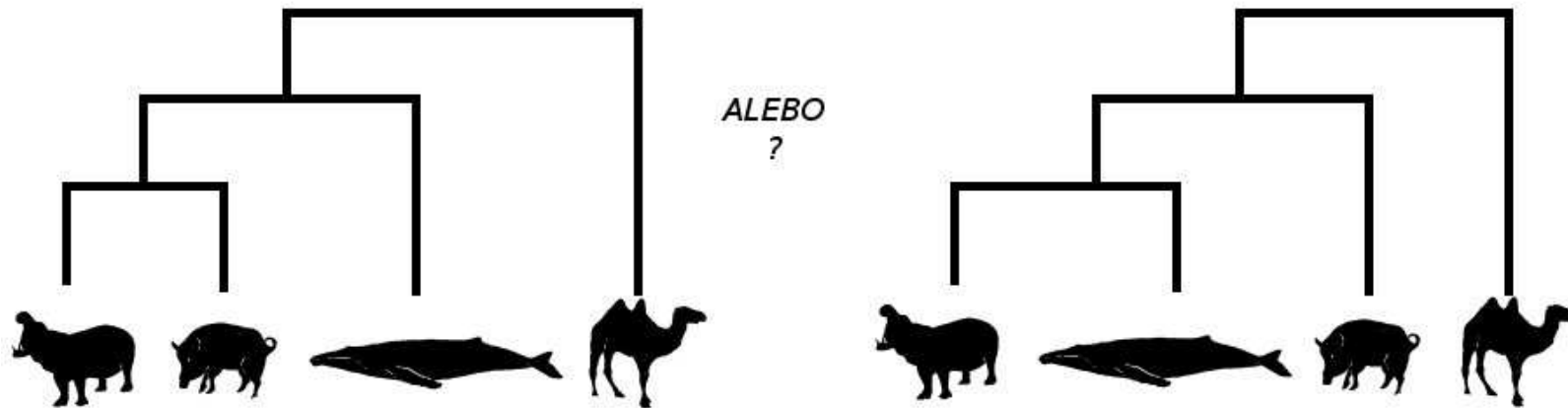
- Domáca úloha 1: odovzdať do budúceho štvrtka 3.11. 9:00 pod dvere kancelárie M163
- Pracujte na journal clube (prečítajte si článok, naplánujte si stretnutie)

<http://compbio.fmph.uniba.sk/vyuka/mbi/>

Evoluční modely a stromy

Tomáš Vinař

27.10.2016



Rekonštrukcia fylogenetických stromov

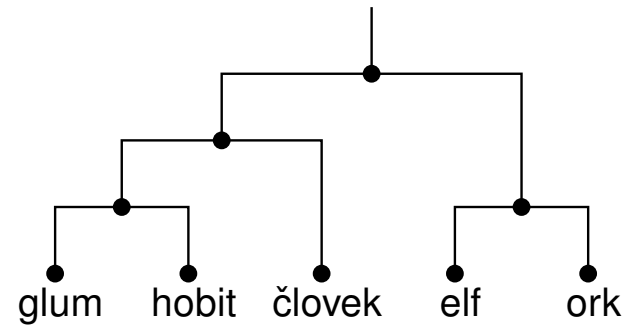
Vstup:

m zarovnaných sekvencií,
každá dĺžky n

človek	C	A	G	T	T	A
elf	A	A	T	A	G	A
Glum	C	C	G	A	G	A
hobit	C	C	G	T	T	C
ork	A	A	T	T	T	A

Výstup:

strom predstavujúci
ich evolučnú históriu

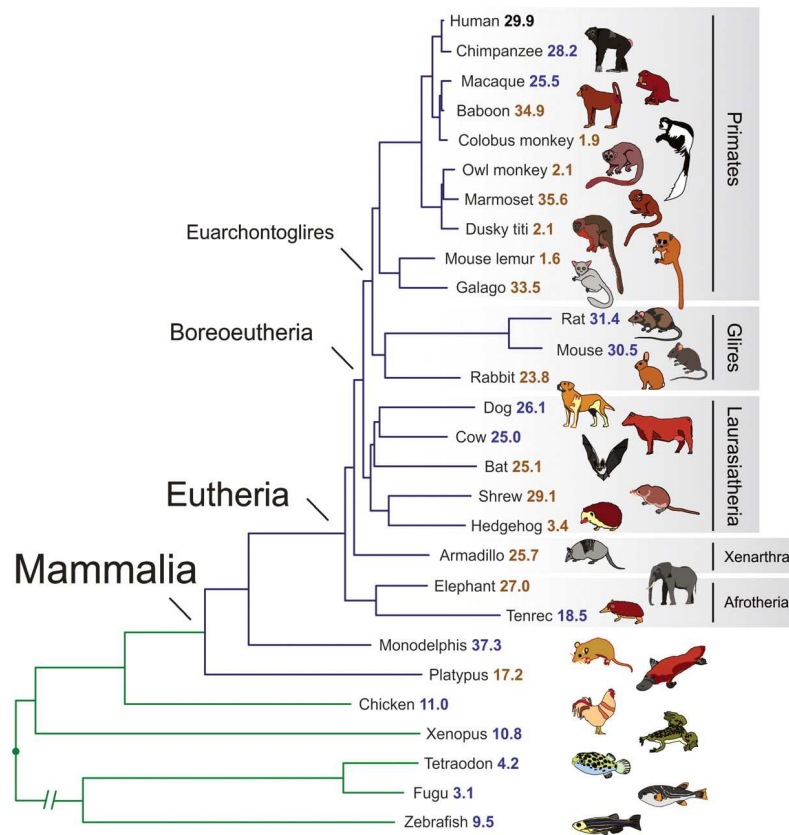


Newick format:

```
((glum,hobit),človek),(elf,ork))
```

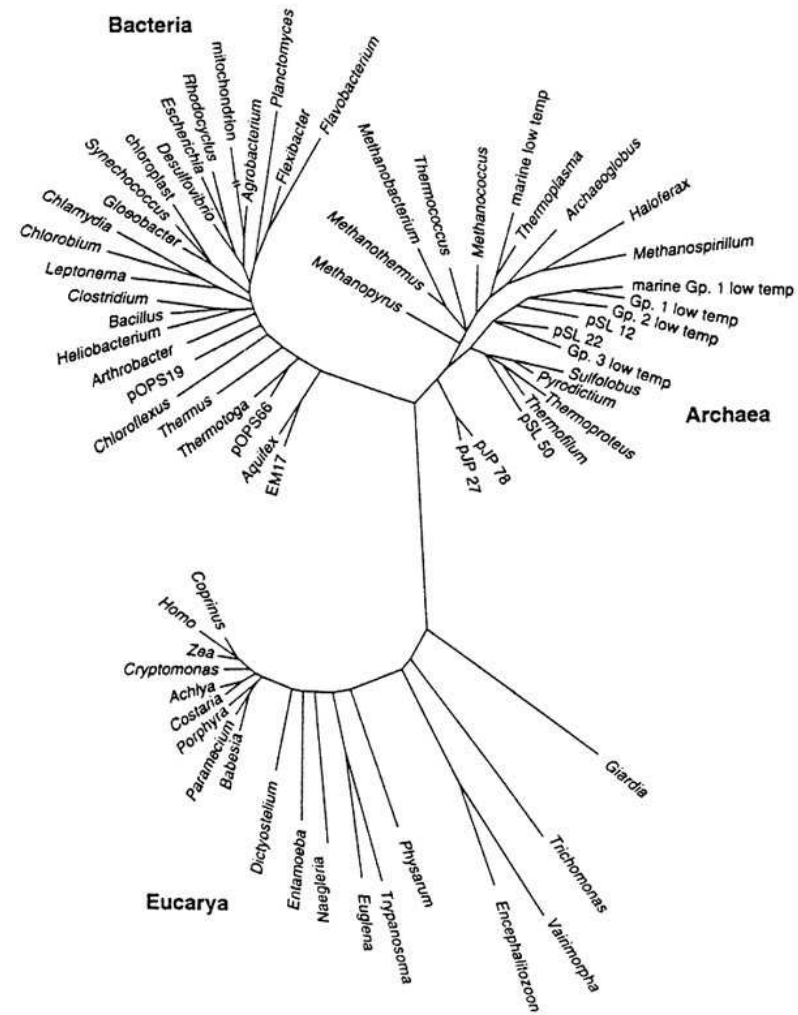
Zakorenené a nezakorenené stromy

[Margulies et al., 2007]



zakorenený pomocou
“outgroup”

[Pace, 1997]



Maximum parsimony (úsporné stromy)

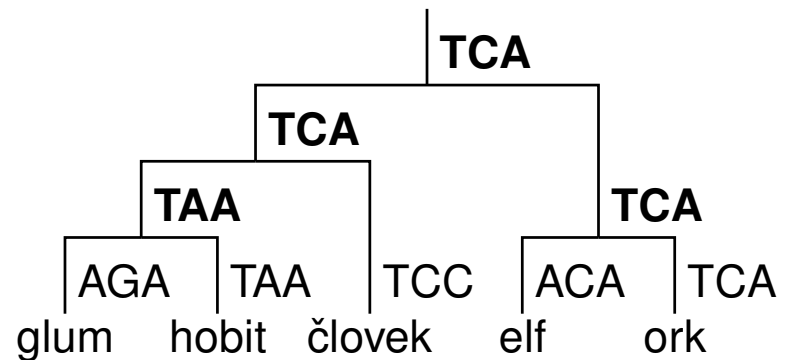
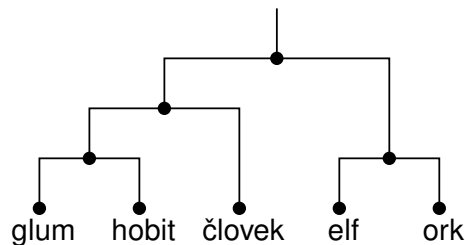
Úloha: Dané sú zarovnané sekvencie súčasných organizmov.

Chceme nájsť fylogenetický strom, ktorý vyžaduje **minimálny počet evolučných zmien**.

Evolučná zmena = mutácia jednej bázy na inú bázu

Podotázka: Pre daný fylogenetický strom, doplniť **ancestrálne sekvencie** tak, aby bol potrebný najmenší počet zmien.

glum	AGA
hobit	TAA
človek	TCC
elf	ACA
ork	TCA



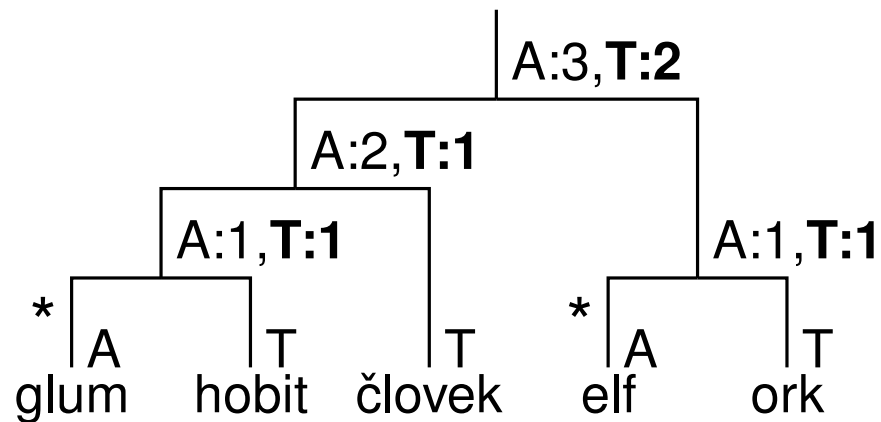
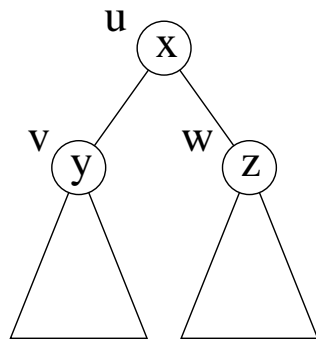
Výpočet ceny konkrétného stromu

Môžeme rátať **dynamickým programovaním** pre každý stĺpec zarovnania zvlášť.

Pre každý vnútorný vrchol u a symbol x :

$N_{u,x}$: koľko zmien treba v podstrome pod u , ak v u bude symbol x ?

$$N_{u,x} = \min_y \{N_{v,y} + [x \neq y]\} + \min_z \{N_{w,z} + [x \neq z]\}$$



Časová zložitosť: $O(m)$, lineárna

Hľadanie najúspornejšieho stromu

NP-ťažký problém

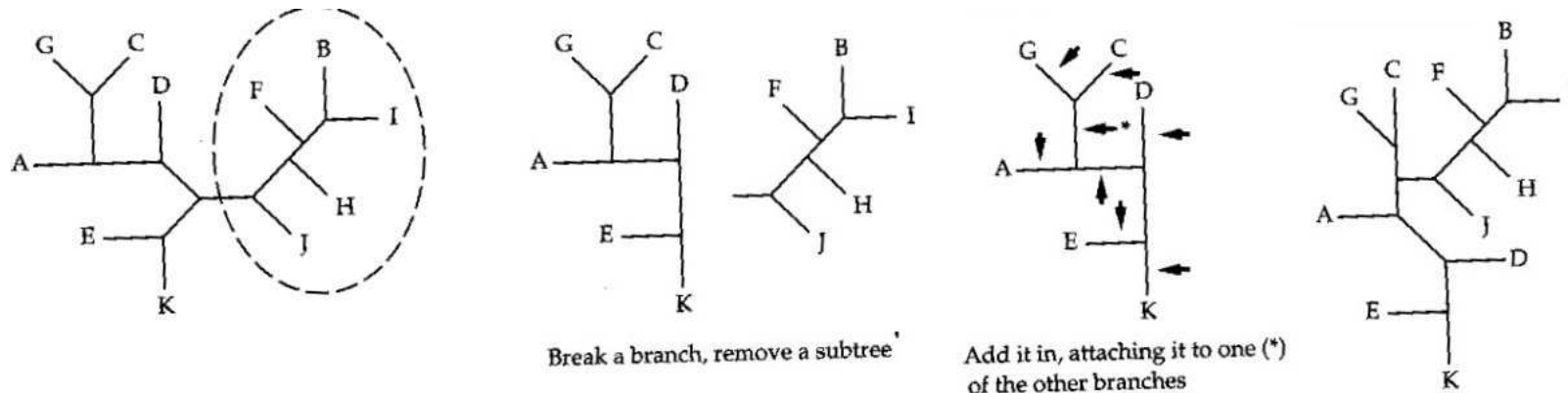
Triviálny algoritmus: vyskúšaj všetky možné stromy.

Pre m druhov $1 \cdot 3 \cdot 5 \cdots (2m - 5) = (2m - 5)!!$

Napr. pre 10 druhov cca 2 milióny, pre 20 druhov $2 \cdot 10^{20}$

Heuristické prehľadávanie:

- Začneme s “rozumným” stromom
- Pomocou stanovených operácií prehľadávame “podobné” stromy;
napr. “subtree pruning and regraft”:



Neighbor Joining (Metóda spájania susedov)

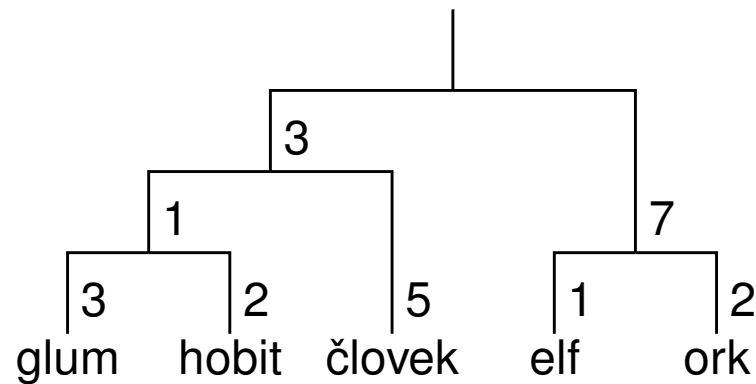
- Nevyužívame detaily rozdielov medzi sekvenciami
- Zosumarizujeme ich pomocou **matice vzdialeností** (D_{ij})

Jednoduchý príklad:

človek	C	A	G	T	T	A		Č	E	G	H	O
elf	A	A	T	A	G	A	človek	0	4	3	2	2
Glum	C	C	G	A	G	A	elf	4	0	3	6	2
hobit	C	C	G	T	T	C	Glum	3	3	0	3	5
ork	A	A	T	T	T	A	hobit	2	6	3	0	4
							ork	2	2	5	4	0

Idea spájania susedov

- Predpokladáme, že vzdialenosti $D_{i,j}$ skutočne zodpovedajú vzdialenostiam v strome (**aditivita**)



$$D_{\text{hobit}, \text{človek}} = 2 + 1 + 5 = 8$$

	glum	hobit	človek	elf	ork
glum	0	5	9	15	16
hobit	5	0	8	14	15
človek	9	8	0	16	17
elf	15	14	16	0	3
ork	16	15	17	3	0

Idea spájania susedov

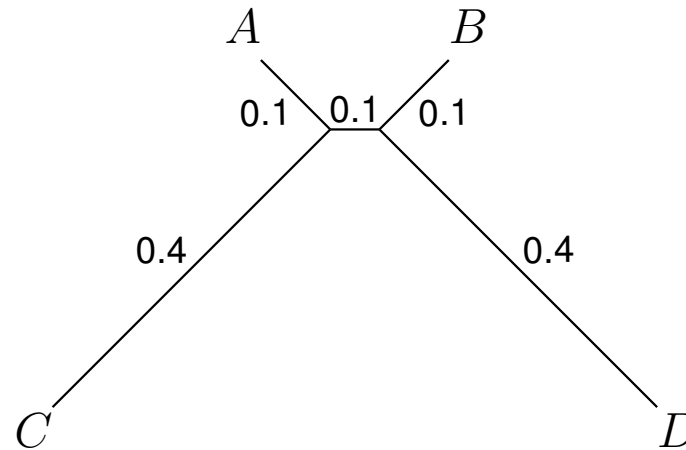
- Predpokladáme, že vzdialenosti $D_{i,j}$ skutočne zodpovedajú vzdialenostiam v strome (**aditivita**)
- Nájďme dva listy i a j , o ktorých vieme **s určitosťou povedať**, že majú vo výslednom strome spoločného otca
- i a j spojíme a nahradíme ich ich otcom k s novými vzdialenosťami:

$$D_{k,\ell} = \frac{D_{i,\ell} + D_{j,\ell} - D_{i,j}}{2}$$

Časová zložitosť: $O(m^3)$

Ako určiť dva listy na spájanie?

(Prečo nie dva najbližšie?)



Vyber listy i, j , ktoré **minimalizujú** nasledujúci výraz:

$$L_{i,j} = (m - 2)D_{i,j} - \underbrace{\sum_{k \neq i} D_{i,k}}_{r_i} - \underbrace{\sum_{k \neq j} D_{j,k}}_{r_j}$$

Spájanie susedov: zhrnutie

- Ak je vstupná matica aditívna a zodpovedá skutočným evolučným vzdialenostiam, spájanie susedov nám dá správny strom
- Čím dlhšie sekvencie, tým spoľahlivejší odhad vzdialenosti a tým väčšia šanca dostať správny strom
- Ako však prejdeme od sekvencií k odhadu vzdialenosti?
Len počítanie rozdielov nestačí

človek	C	A	G	T	T	A		Č	E	G	H	O
elf	A	A	T	A	G	A	človek	0	4	3	2	2
Glum	C	C	G	A	G	A	elf	4	0	3	6	2
hobit	C	C	G	T	T	C	Glum	3	3	0	3	5
ork	A	A	T	T	T	A	hobit	2	6	3	0	4
							ork	2	2	5	4	0

Problém so vzdialenosťami

- Počas evolúcie sa môže stať, že tá istá báza zmutuje **viackrát** (trebárs aj späť na originálnu bázu)
- Pri počítaní rozdielov ale vidíme nanajvýš jednu zmenu na každej pozícii \Rightarrow odhad vzdialenosti menší ako v skutočnosti
- Chceme korekciu na odhadovaný počet mutácií, ktoré sa naozaj stali

Jukesov-Cantorov model evolúcie

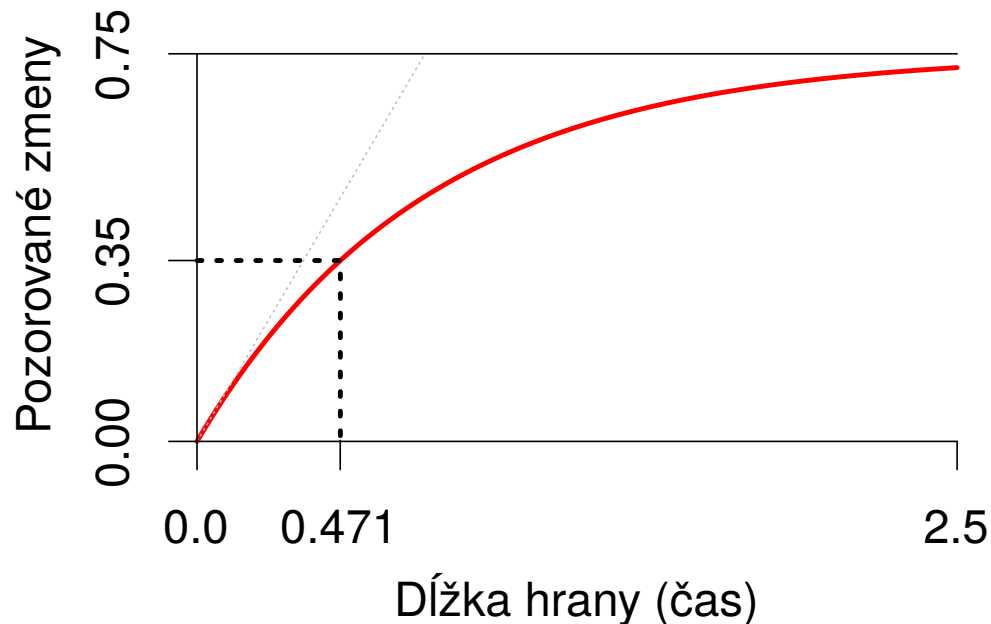
Pravdepodobnosť zmeny bázy na inú:

$$\Pr(X_{t_0+t} = C \mid X_{t_0} = A) = \frac{1}{4}(1 - e^{-\frac{4}{3}\alpha t})$$

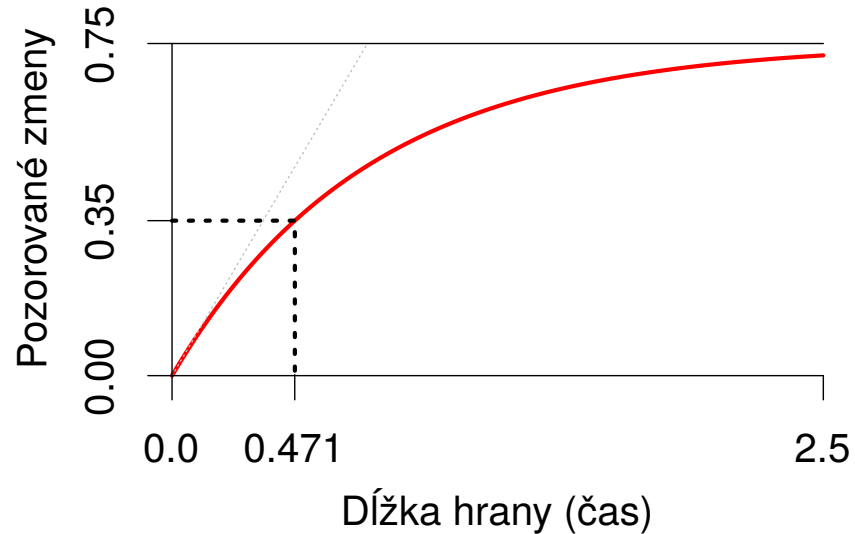
α : rýchlosť evolúcie (počet substitúcií na jednotku času)

Očakávaný počet pozorovaných zmien na bázu za čas t :

$$D(t) = \Pr(X_{t_0+t} \neq X_{t_0}) = \frac{3}{4}(1 - e^{-\frac{4}{3}\alpha t})$$



Späť ku spájaniu susedov (Neighbor Joining)



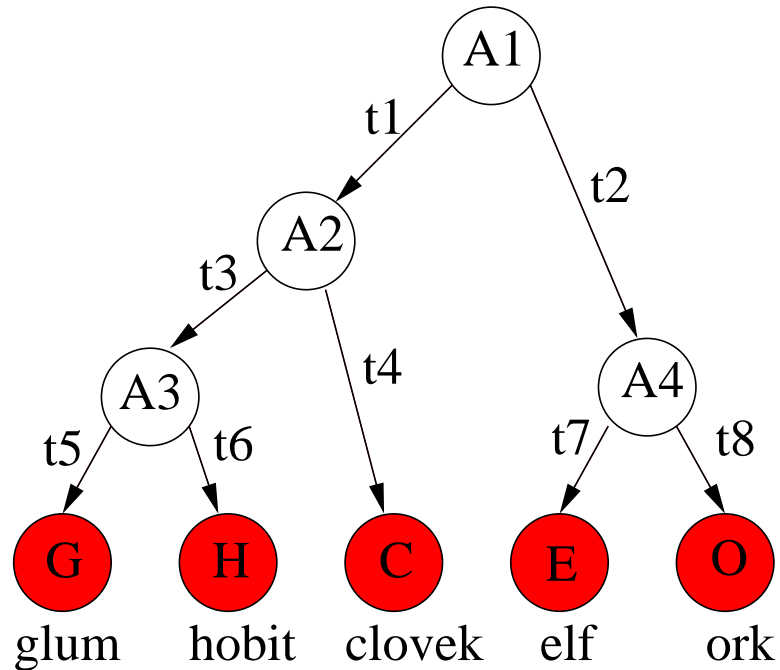
- Podľa takéhoto modelu môžeme korigovať pozorované vzdialenosti

$$D = \frac{3}{4} (1 - e^{-\frac{4}{3}\alpha t}) \quad \Rightarrow \quad \alpha t = -\frac{3}{4} \ln\left(1 - \frac{4}{3}D\right)$$

- Často lepšie zložitejšie modely, ktoré zahŕňajú rôzne frekvencie báz, pomer tranzícií a transverzií, variabilnú rýchlosť evolúcie na rôznych pozíciách (pozri [Felsenstein, 2004, kap.13])

Najvierohodnejšie stromy (Maximum likelihood)

- Strom môžeme chápať ako **jednoduchý generatívny model**



- $\Pr(G, H, C, E, O, A1, \dots, A4) = \Pr(A1) \cdot \Pr(A2 | A1, t_1) \cdot \Pr(A4 | A1, t_2) \cdot \Pr(A3 | A2, t_3) \cdot \Pr(G | A3, t_5) \cdot \Pr(H | A3, t_6) \cdot \Pr(C | A2, t_4) \cdot \Pr(E | A4, t_7) \cdot \Pr(O | A4, t_8)$

- $\Pr(G, H, C, E, O, A_1, \dots, A_4) = \Pr(A_1) \cdot \Pr(A_2 | A_1, t_1) \cdot \Pr(A_4 | A_1, t_2) \cdot \Pr(A_3 | A_2, t_3) \cdot \Pr(G | A_3, t_5) \cdot \Pr(H | A_3, t_6) \cdot \Pr(C | A_2, t_4) \cdot \Pr(E | A_4, t_7) \cdot \Pr(O | A_4, t_8)$
- Pre **daný strom** a **dané dĺžky hrán** možno jednotlivé pravdepodobnosti spočítať použitím evolučného modelu (napr. Jukes-Cantor)

- **Vierohodnosť (likelihood) stromu:**

$$\Pr(G, H, C, E, O) = \sum_{A_1, \dots, A_4} \Pr(G, H, C, E, O, A_1, \dots, A_4)$$

- Rátame pomocou **Felsensteinovho algoritmu** (jednoduché dynamické programovanie, podobne ako pre parsimony)
- \Rightarrow Pre daný strom a dĺžky hrán vieme spočítať vierohodnosť v čase $O(m)$

Ako nájsť najvieryhodnejší strom?

- Problém je NP-ťažký ;
navyše komplikovaný tým, že na výpočet vierohodnosti **potrebujeme aj dĺžky hrán**
- Opäť použijeme heuristické vyhľadávanie:
 - Začneme s “rozumným” stromom
 - Vypočítame vierohodnosť tohto stromu:
 - * Začneme s “rozumnými” dĺžkami hrán
 - * Vypočítame vierohodnosť stromu s dĺžkami
 - * Mierne zmeníme dĺžky tak, aby sa zlepšila vierohodnosť a opakujeme
 - Pomocou stanovených operácií (ako v prípade parsimony) skúsime “podobné” stromy, až kým nevieme zlepšiť

“Správnosť” fylogenetických algoritmov: Konzistentnosť

- “Rozumne” správajúce sa algoritmy: ak množstvo dát (n) rastie, ich odpoveď by sa mala približovať ku správnej odpovedi.
- Hovoríme, že algoritmus pre hľadanie fylogenetického stromu je **konzistentný**, ak v prípade, že n ide do nekonečna, pravdepodobnosť správneho stromu konverguje k 1.

Porovnanie algoritmov

	Zložitosť	Konzistentný	Využitie dát
Parsimony (úspornosť)	NP-ťažký	NIE	celé sekvencie
Neighbor Joining	$O(m^3)$	ÁNO	iba vzdialenosti
Likelihood (vierohodnosť)	NP-ťažký	ÁNO	celé sekvencie

Odkiaľ zohnať dáta pre fylogenetiku?

- **Mitochondriálna DNA (mtDNA):**

- Krátky cirkulárny genóm uložený v mitochondriách (človek: cca 16KB)
- Dedí sa po materskej línii (žiadna rekombinácia)
- Rýchlejšie mutácie – vhodný nielen pre druhy, ale aj jedince
- Ľahko sa sekvenuje; osekvenovaný pre mnoho organizmov

- **Ribozomálna RNA (rRNA):**

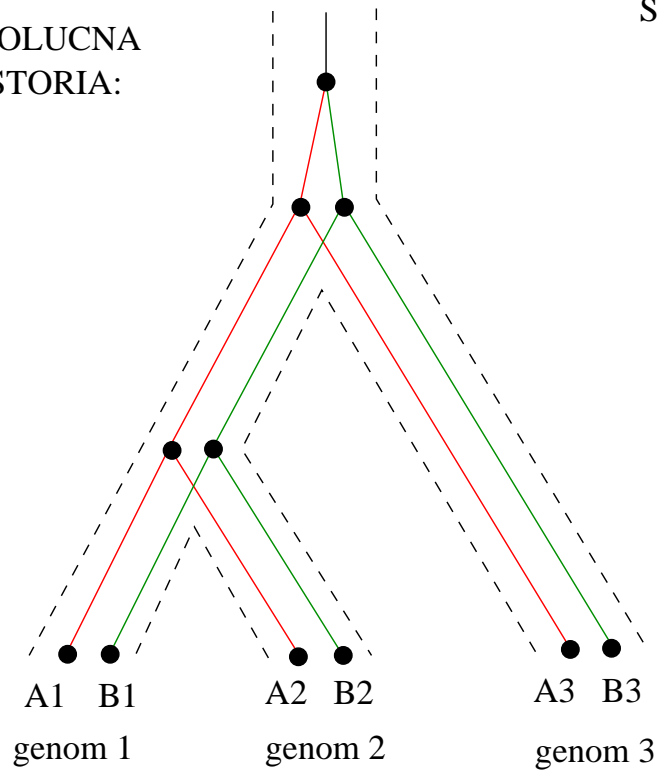
- Nepostrádateľná pri syntéze proteínov v ribozómoch
- ⇒ veľmi dobre zachovaná aj medzi druhmi
- RDPII databáza

- **DNA sekvencie:** Čo tak:

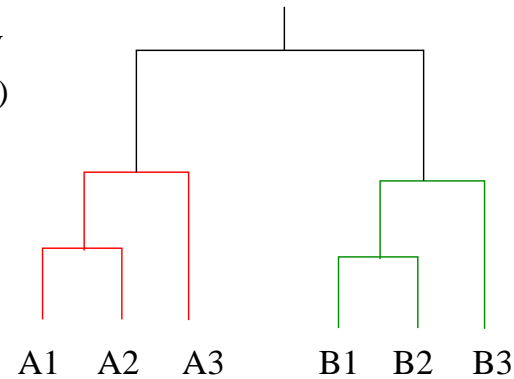
- Vybrať si sympatický gén
- Nájsť jeho homológy v iných genómoch
- Použiť tieto na konštrukciu fylogenetického stromu

Problém!!! Duplikácia génov (a vo všeobecnosti DNA duplikácia)

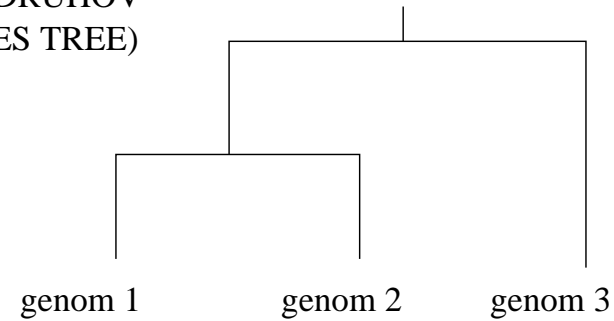
EVOLUCNA
HISTORIA:



STROM GENOV
(GENE TREE)

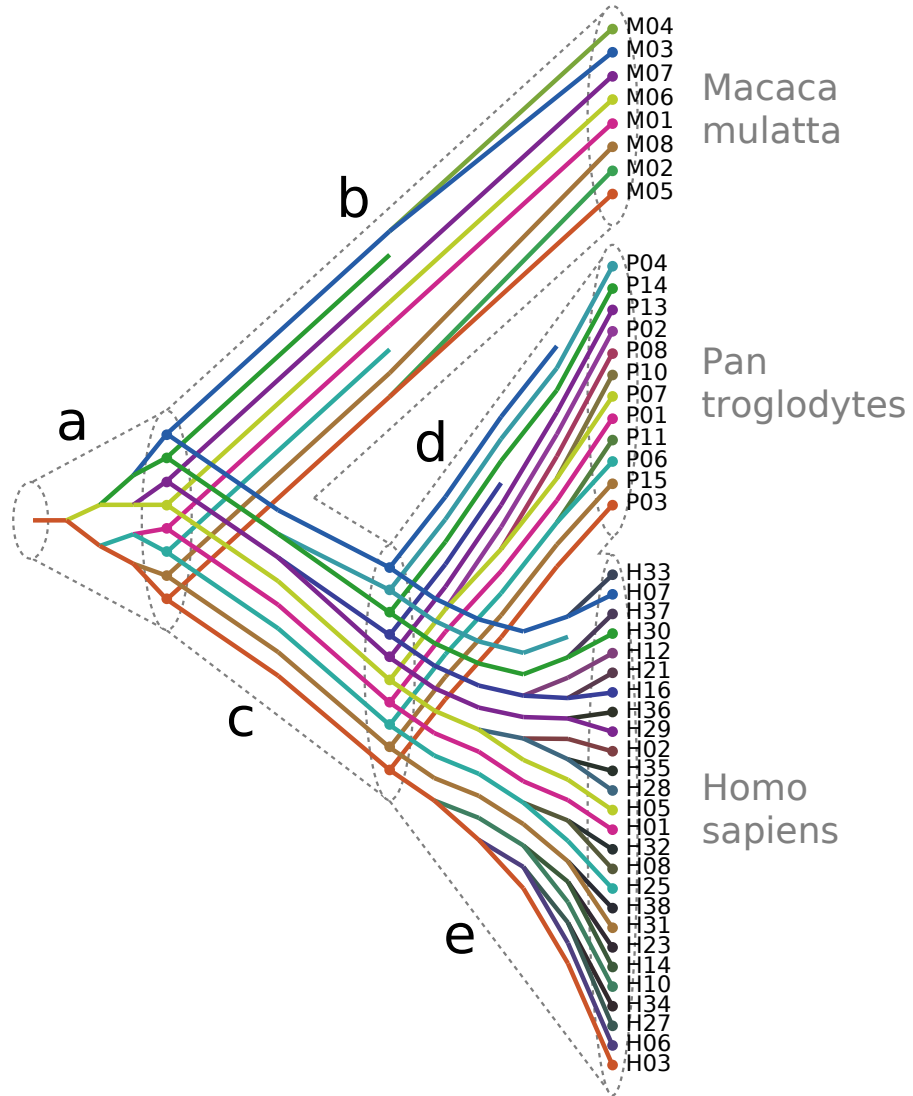


STROM DRUHOV
(SPECIES TREE)



- **Homológ:** vyvinuli sa zo spoločného predka, podobná sekvencia
- **Ortológ:** najbližší spoločný predok je speciácia (napr. A1/A3)
- **Paralóg:** najbližší spoločný predok je duplikácia (napr. A1/B1, A1/B2)

Zložitejší příklad:



Zhrnutie:

- Modely evolúcie nukleotidov nám dávajú možnosť:
 - Odhadovať skutočnú evolučnú vzdialenosť (počet substitúcií) z počtu pozorovaných zmien medzi sekvenciami
 - Počítať pravdepodobnosť, že uvidíme zmenu nukleotidu za určitý čas t
- Tri metódy na vytváranie evolučných stromov:
 - Úsporné stromy (parsimony)
 - Spájanie susedov (neighbour joining)
 - Vierohodnosť stromov (maximum likelihood)
- Génové a druhové stromy; komplikácie pri vytváraní stromov

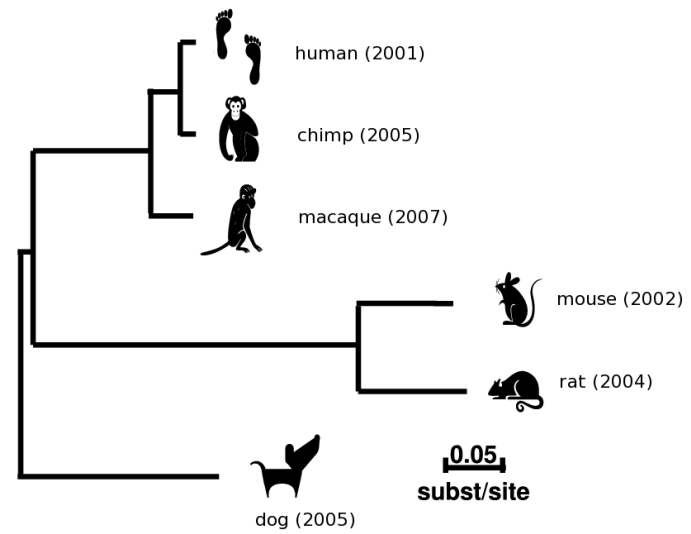
Organizačné poznámky

- **Domáca úloha 2**
bude na stránke budúci týždeň
- Nezabudnite na prvé stretnutie ohľadom journal clubu!!!

Komparatīvna genomika

Tomás Vinar̄

3.11.2016



Komparatívna genomika

- Štúdium evolúcie genómov
 - Mutácie jednotlivých báz DNA (táto prednáška)
 - Krátke inzercie a delécie
 - Väčšie udalosti: prestavby genómu, duplikácie
- Typy mutácií:
 - Neutrálne
 - Škodlivé (deleterious)
 - ⇒ **Purifikačný výber (purifying selection)**
 - Prospešné (advantageous)
 - ⇒ **Pozitívny výber (positive selection)**
- Na základe porovnávania genómov chceme nájsť oblasti s nezvyčajnou evolučnou históriou (zachovávanie dôležitých funkcií, vývoj nových funkcií)

Komparatívna genomika

- Zostavíme viacnásobné zarovnanie genómov
(zarovnané miesta by mali pochádzať z tej istej sekvencie spoločného predka)

```
Human  AGTGGCTGCCAGGCTG---GGATGCTGAGGCCTTGTTTGCAGGGAGGT
Rhesus AGTGGCTGCCAGGCTG---GGTTGCTGAGGCCTTGTTTGCCGGGAGGT
Mouse  GGTGGCTGCCGGGCTG---GGTGGCTGAGGCCTTGTTGGTGGGGTGGT
Dog     AGTGGCTGCCCGGCTG---GGTGGCTGAGGCCTTATTTGCAGGGAGGT
Horse  GATGGCTGCCGGGCTG---GGCTGCCGAGGCCTTGTTTCGTGGGGAGGT
Armadillo AGTGGCTGCCGGGCTG---GGAGGCCAAGGCCTTGTTTCGCGGGCAGGT
Chicken AGTGGCTGCCAGTCTGCGCCGTGGCCGACGTCTTGCTCGGGGAAGGT
X. tropicalis AATGGCTTCCATTTTGTGCCGCTGCTGAGGTCTTGTTCTGGGGAAGAT
```

- **Metódy:** Kombinujeme techniky na anotáciu (HMM) a pravdepodobnostné modely evolúcie

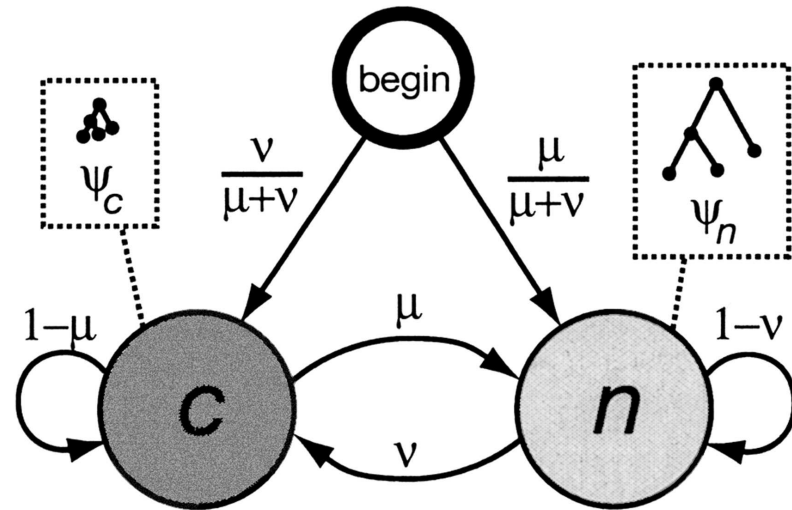
Príklad 1: Hľadanie funkčných oblastí sekvencií

Dôsledky purifikačného výberu:

- Funkčné časti sekvencie zostávajú zachované, menia sa pomalšie
- Nefunkčné sekvencie sa vyvíjajú rýchlejšim tempom
- **Príklad:** gény kódujúce proteíny, porovnanie človek myš
 - kódujúce časti: 85% zhoda (zarovnanie na 98% dĺžky)
 - intróny: 69% zhoda (zarovnanie na 48% dĺžky)
- **Úloha:** Hľadáme **nadmerne dobre zachované sekvencie**
- Veľká časť bude zodpovedať známym funkčným elementom (kódujúce gény, regulačné regióny, a pod.)
- Zachované sekvencie ktoré sa neprekrývajú so známymi funkčnými elementami — zaujímavé objekty pre výskum

PhastCons: detekcia dobre zachovaných sekvencií

Fylogenetické HMM: kombinácia HMM a fylogenetického stromu.



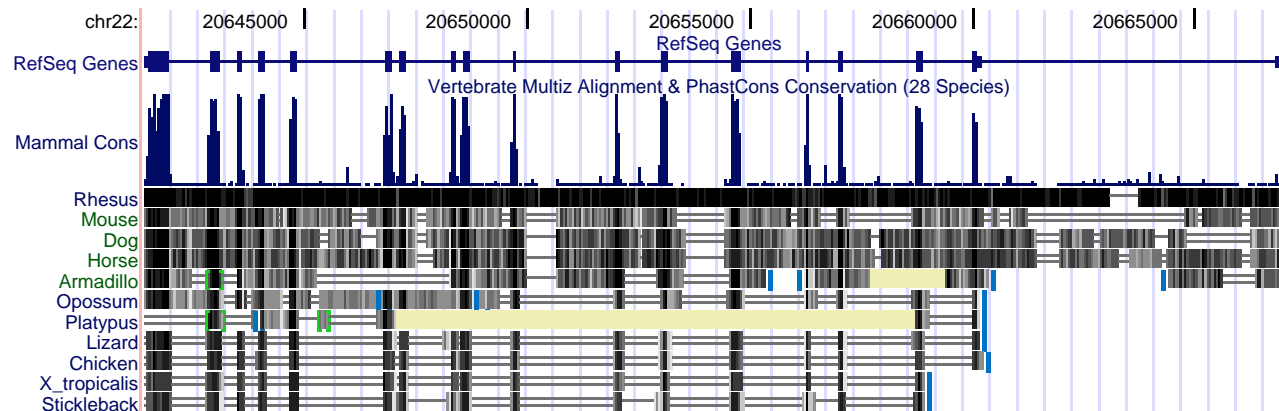
x =

TCGCGACATATACGA	...
TTGGGGCATGTGGGT	...
AGCAGACGTCCGCAA	...

- Dva stavy: zachovaná sekv., neutrálna sekv.
- V každom stave generujeme celý stĺpec zarovnania
- Zachovaná sekvencia má kratšie hrany stromu

Použitie fylogenetického HMM

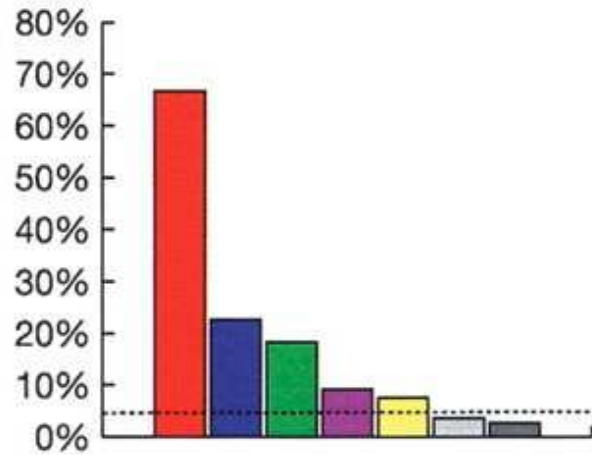
- Model určuje rozdelenie pravdepodobnosti cez zarovnanie a anotácie
(tu: anotácia = označenie zachovaných sekvencií)
- Pre dané zarovnanie hľadáme najpravdepodobnejšiu anotáciu
- Kombinácia Viterbiho a Felsensteinovho algoritmu



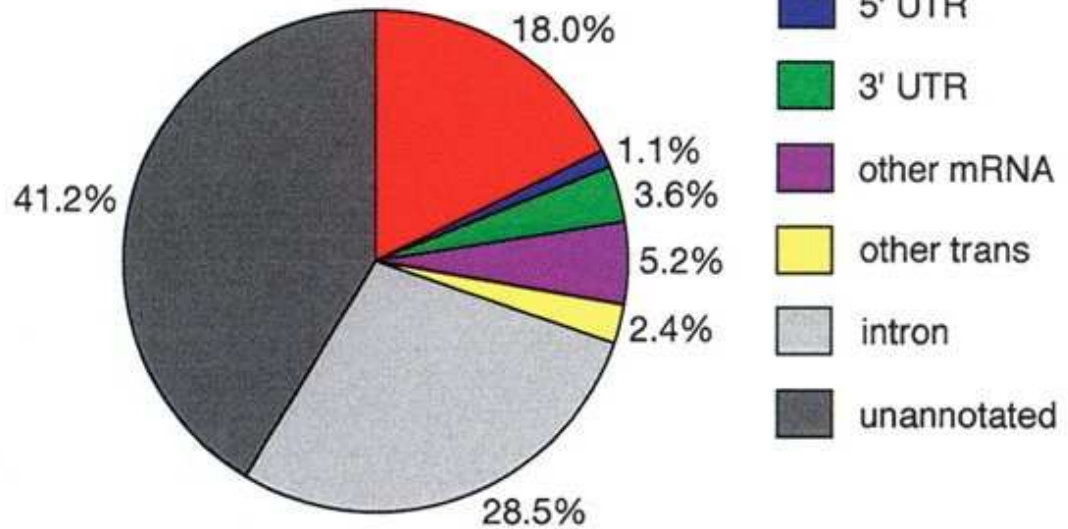
Výsledky celogenómovej aplikácie PhastCons-u

Zarovnania genómov človeka, myši, sliepky, fugu

Coverage of Annotation Types by Conserved Elements



Composition of Conserved Elements by Annotation Type



Genetický kód

Alanine (A)

GC*

Cysteine (C)

TGC

TGT

Aspartic acid (D)

GAC

GAT

Glutamic acid (E)

GAA

GAG

Phenylalanine (F)

TTC

TTT

Glycine (G)

GG*

Histidine (H)

CAC

CAT

Isoleucine (I)

ATA

ATC

ATT

Lysine (K)

AAA

AAG

Leucine (L)

CT*

TTA

TTG

Methionine (M)

ATG

Asparagine (N)

AAC

AAT

Proline (P)

CC*

Glutamine (Q)

CAA

CAG

Arginine (R)

CG*

AGA

AGG

Serine (S)

TC*

AGT

AGC

Threonine (T)

AC*

Valine (V)

GT*

Tryptophan (W)

TGG

Tyrosine (Y)

TAC

TAT

Stop codon (*)

TAA

TAG

TGA

Fylogenetické HMM pre hľadanie génov

- Použijeme stavy z hľadača génov
- Pre každý stav máme evolučný model (maticu rýchlostí, dĺžky hrán)
- Trojperiodickosť frekvencií mutácií pomáha nájsť gény

Ako veľmi pomôžu zarovnaná zlepšiť presnosť

Program	Exóny		Gény	
	sn	sp	sn	sp
AUGUSTUS (1 genóm)	52%	63%	24%	17%
NSCAN (zarovnanie)	68%	82%	35%	37%

Guigo et al 2006, evaluácia na 1% ľudského genómu

Príklad 2: Hľadanie génov pod vplyvom pozitívneho výberu

- **Pozitívny výber** = proces, ktorým sa v genóme ustália **prospešné mutácie**
- Neobvykle vysoké množstvo mutácií, ktoré by mohli súvisieť so zmenou funkcie
- V rámci génov, ktoré kódujú proteíny:
 - **Synonymné mutácie** nemenia zakódovanú aminokyselinu
napr. ACA (Thr) \Rightarrow ACT (Thr)
 - **Nesynonymné mutácie** menia zakódovanú aminokyselinu
napr. ACA (Thr) \Rightarrow AAA (Lys)
- Vytvoríme pravdepodobnostný model evolúcie, ktorý bude rozlišovať synonymné a nesynonymné mutácie \Rightarrow identifikácia sekvencií s neobvykle vysokým podielom nesynonymných mutácií

Od Jukes-Cantorovho modelu ku všeobecnejším modelom mutácií

- Jukes-Cantor predpokladá, že každá mutácia rovnako pravdepodobná
- Vo všeobecnosti zavedieme μ_{xy} – **rýchlosť substitúcie** z bázy x na bázu y
- Matica rýchlostí (substitution rate matrix)

$$\begin{pmatrix} -\mu_A & \mu_{AC} & \mu_{AG} & \mu_{AT} \\ \mu_{CA} & -\mu_C & \mu_{CG} & \mu_{CT} \\ \mu_{GA} & \mu_{GC} & -\mu_G & \mu_{GT} \\ \mu_{TA} & \mu_{TC} & \mu_{TG} & -\mu_T \end{pmatrix}$$

$$\begin{pmatrix} -\mu_A & \mu_{AC} & \mu_{AG} & \mu_{AT} \\ \mu_{CA} & -\mu_C & \mu_{CG} & \mu_{CT} \\ \mu_{GA} & \mu_{GC} & -\mu_G & \mu_{GT} \\ \mu_{TA} & \mu_{TC} & \mu_{TG} & -\mu_T \end{pmatrix}$$

- **Rovnovážny stav:** frekvencie $\pi_A, \pi_C, \pi_G, \pi_T$ nemení sa v čase
- Pre daný čas t , môžeme vypočítať pravdepodobnosť každej substitúcie (**transition probabilities**):

$$\Pr(X = C \mid Y = A, t)$$

Znižovanie počtu parametrov — HKY matica

Hasegawa, Kishino a Yano

$$\begin{pmatrix} -\mu_A & \beta\pi_C & \alpha\pi_G & \beta\pi_T \\ \beta\pi_A & -\mu_C & \beta\pi_G & \alpha\pi_T \\ \alpha\pi_A & \beta\pi_C & -\mu_G & \beta\pi_T \\ \beta\pi_A & \alpha\pi_C & \beta\pi_G & -\mu_T \end{pmatrix} \quad \mu_{x,y} = \begin{cases} \alpha\pi_y & \text{ak } x \Leftrightarrow y \text{ je tranzícia} \\ \beta\pi_y & \text{ak } x \Leftrightarrow y \text{ je tranzverzia} \end{cases}$$

- **rýchlosť tranzícií (transition rate)** α : $C \Leftrightarrow T, A \Leftrightarrow G$
- **rýchlosť tranzverzií (transversion rate)** β : $\{C, T\} \Leftrightarrow \{A, G\}$
- Máme iba štyri parametre: $\pi_A, \pi_C, \pi_G, \kappa = \alpha/\beta$

Substitučný model pre kodóny

Namiesto jednotlivých báz uvažujeme trojice

Rýchosť zmeny z kodónu i na kodón j :

$$\mu_{i,j} = \begin{cases} 0, & \text{ak sa } i, j \text{ líšia na } > 1 \text{ pozíciách,} \\ \alpha\pi_j, & \text{synonymné tranzície,} \\ \beta\pi_j, & \text{synonymné transverzie,} \\ \omega\alpha\pi_j, & \text{nesynonymné tranzície,} \\ \omega\beta\pi_j, & \text{nesynonymné transverzie.} \end{cases}$$

Príklad: $\mu_{AAC,GGC} = 0$, $\mu_{CTA,CTT} = \beta\pi_{CTT}$,

$\mu_{CTA,CCA} = \omega\alpha\pi_{CCA}$

Parametre: Frekvencie kodónov π_j , ω , $\kappa = \alpha/\beta$

Prirodzený výber: neutrálna evolúcia $\omega = 1$, pozitívny výber $\omega > 1$,
purifikačný výber $\omega < 1$

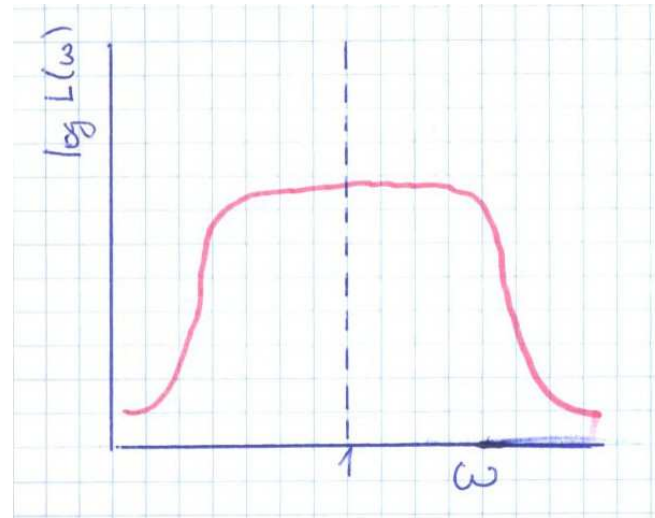
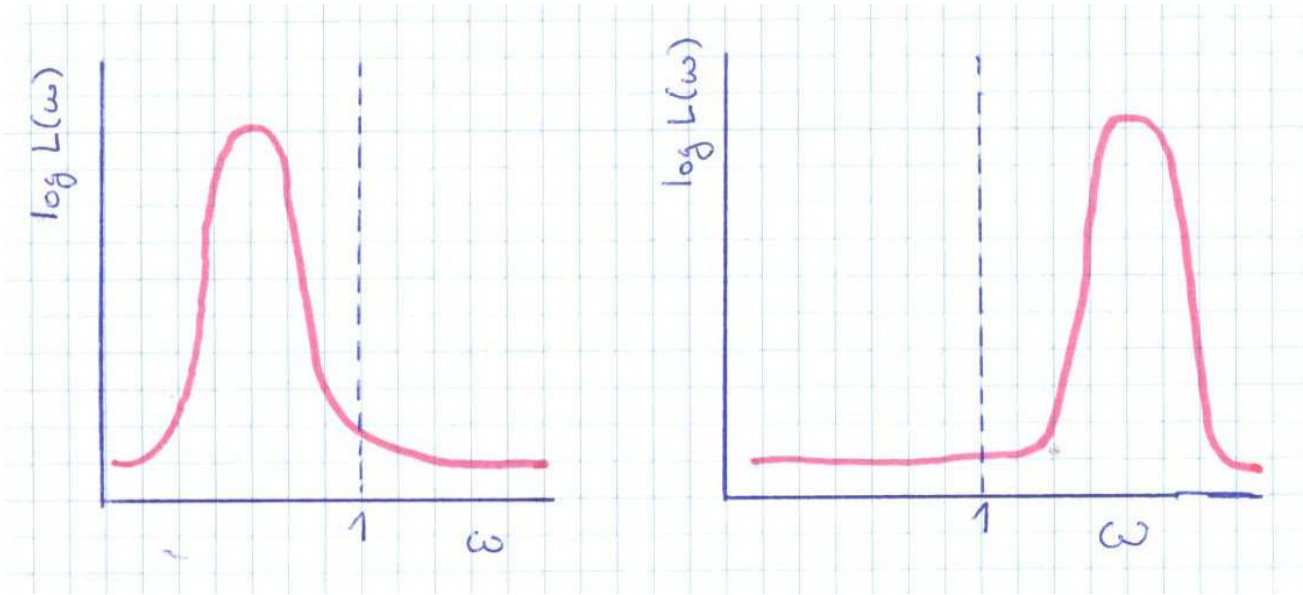
Aplikácia kodónového substitučního modelu

	F	V	I	H	D	S	E	G	D	G	E	C	M	Q	E
človek	TTT	GTG	ATC	CAC	GAC	TCC	GAG	GGG	GAC	GGC	GAG	TGC	ATG	CAG	GAG
kosmák	TTT	GTG	ATC	CAC	GAG	AAC	AAC	AAG	GAC	GGC	GAG	TGC	ATG	CAG	GAT
	F	V	I	H	E	N	N	K	D	G	E	C	M	Q	D

- Na základe celých genómov môžeme odhadnúť základné parametre modelu $\pi_A, \pi_C, \pi_G, \pi_T, \kappa$
- Pre dané ω a t vieme spočítať vierohodnosť

$$L(\omega, t) = \Pr(C, K \mid \omega, t)$$

- Sledujeme, ako sa mení $L(\omega) = \max_t L(\omega, t)$ pre rôzne hodnoty ω



Test pomerov vierohodností (Likelihood-ratio test)

- $L(\omega)$ môže byť najväčšie pre $\omega > 1$,
ale môže to byť spôsobené len štatistickou variáciou v dátach
 \Rightarrow potrebujeme štatistický test
- Spočítame vierohodnosť $L_A = \max_{\omega < 1} L(\omega)$
- Spočítame vierohodnosť $L_B = \max_{\omega} L(\omega)$ (bez obmedzenia ω)
- Vždy platí $L_B \geq L_A$
- Ak skutočné $\omega < 1$, $L_A \approx L_B$ (nulová hypotéza)
 nás zaujímajú prípady $L_B \gg L_A$
 \Rightarrow gén pod vplyvom pozitívneho výberu (alt. hypotéza)

Za predpokladu, že $\omega < 1$, platí $2 \log(L_B/L_A) \approx \chi_1^2$
 \Rightarrow možno priradiť P-hodnotu nulovej hypotéze $\omega < 1$

Hľadanie génov pod vplyvom pozitívneho výberu: Zhrnutie

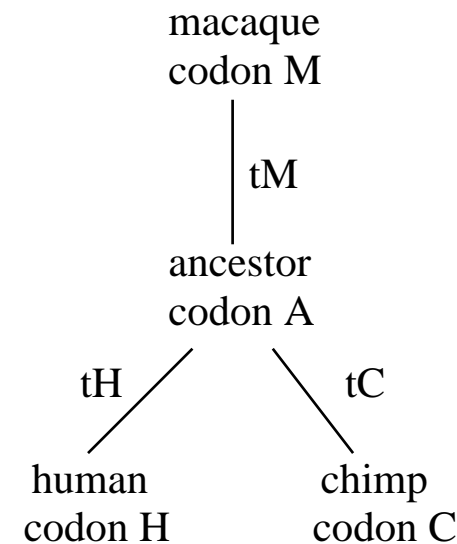
- Nájdem zariadenie toho istého génu z dvoch organizmov (na úrovni kodónov)
- Odhadneme základné parametre kodónového modelu na základe porovnania celých genómov
- Parameter ω modeluje selekciu
- Spočítame vierohodnosť $L_A = \max_{\omega < 1} L(\omega)$
a vierohodnosť $L_B = \max_{\omega} L(\omega)$
- Na základe štatistiky $2 \log(L_B/L_A)$ priradíme P-hodnotu nulovej hypotéze $\omega < 1$
- Gény s malou P-hodnotou sú pod vplyvom pozitívneho výberu

“Jednoducho” rozšíriteľné na porovnanie viacerých organizmov

$$\Pr(A, H, C, M \mid \omega, t_H, t_C, t_M) = \pi_A \cdot \Pr(H \mid A, t_H) \cdot \Pr(C \mid A, t_C) \cdot \Pr(M \mid A, t_M)$$

Zbavíme sa ancestrálnych sekvencií:

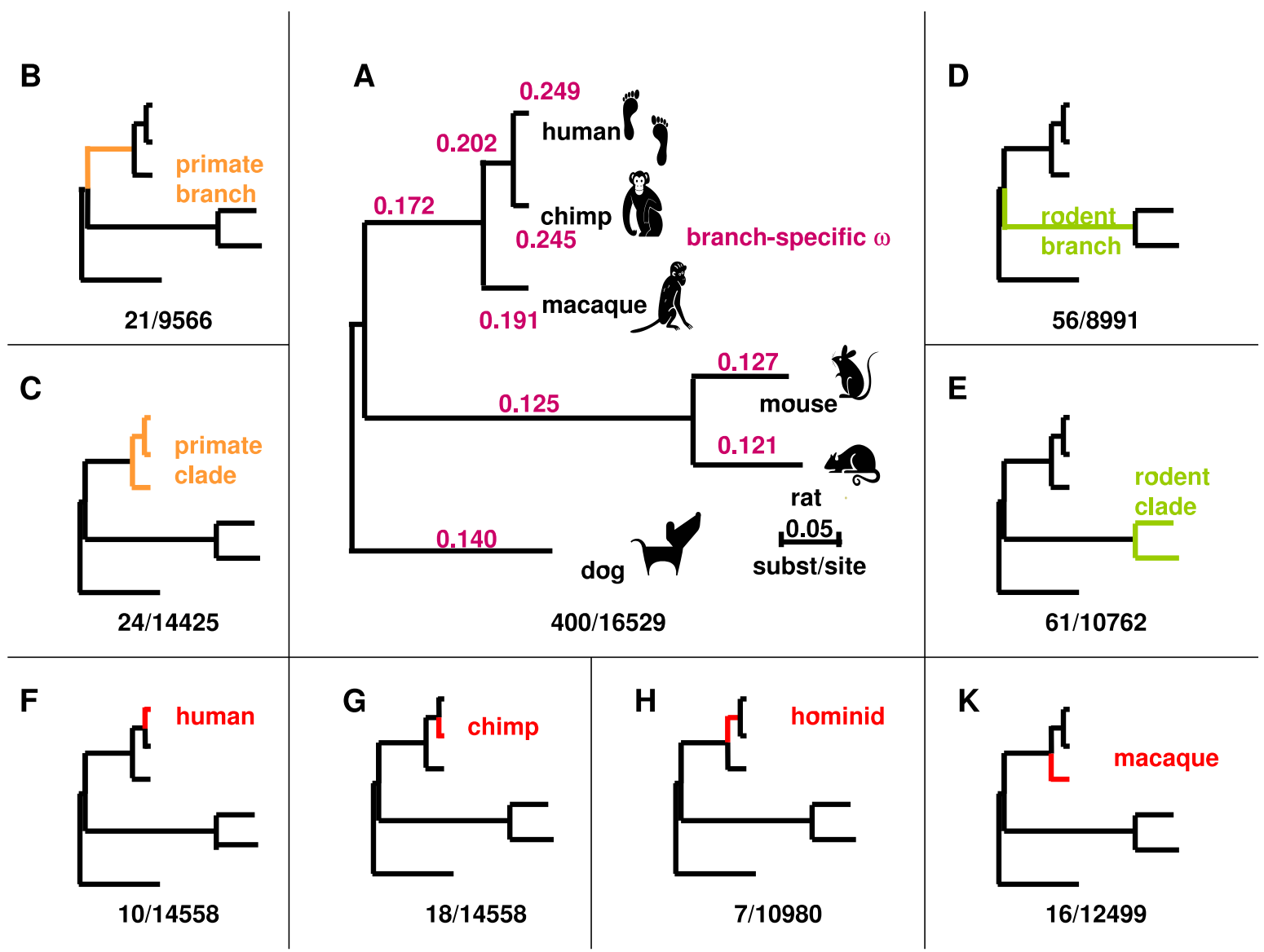
$$\Pr(H, C, M \mid \omega, t_H, t_C, t_M) = \sum_A \Pr(A, H, C, M \mid \omega, t_H, t_C, t_M)$$



Vierohodnosť ω :

$$L(\omega) = \max_{t_H, t_C, t_M} \Pr(H, C, M \mid \omega, t_H, t_C, t_M)$$

- Existuje program PAML, ktorý takúto vierohodnosť počíta
- K dispozícii zložitejšie modely, napr. s meniacim sa ω v rámci génu



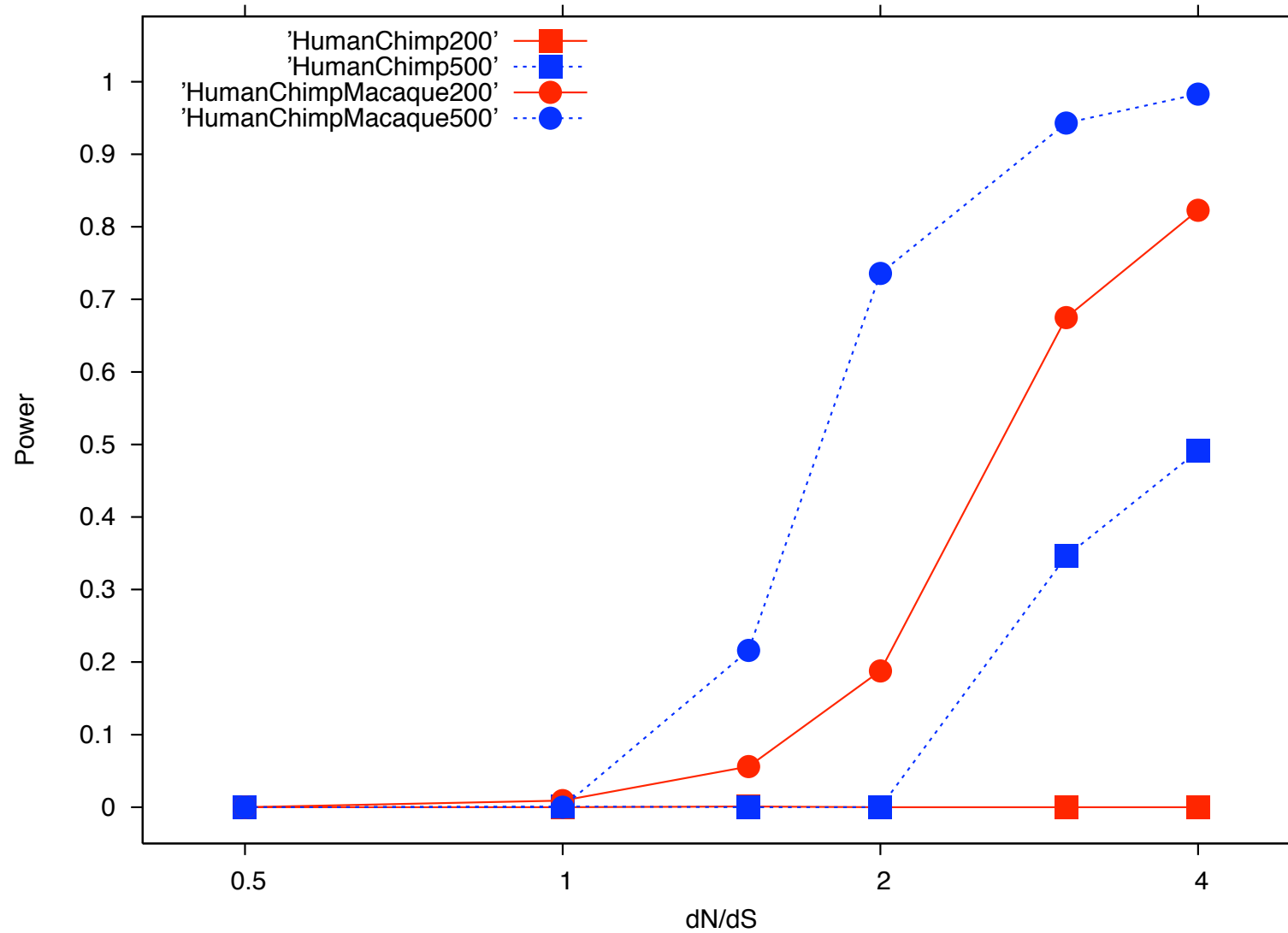
Funkčné kategórie obohatené o gény s pozitívnym výberom

Defense: cellular defense response, antigen processing and presentation, response to virus, response to bacterium

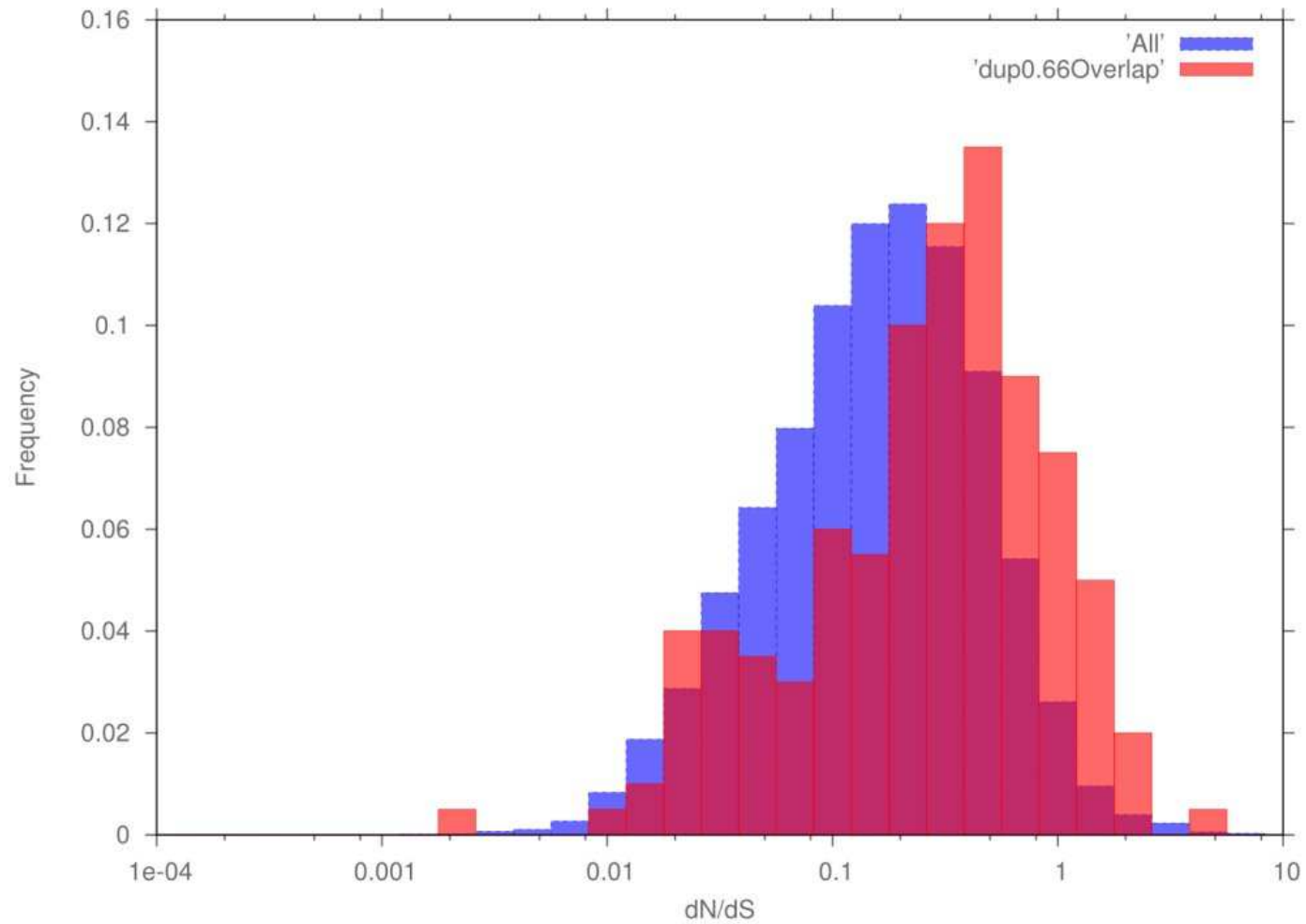
Immunity: adaptive immune response, adaptive immune response somatic recomb, lymphocyte mediated immunity, immunoglobulin mediated immune response, B cell mediated immunity, innate immune response, complement activation alternative pathway, regulation of immune system process, positive regulation of immune response, humoral immune response, complement activation classical pathway, humoral immune response circulating immunoglob, complement activation, activation of plasma proteins mute inflam resp, akute inflammatory response, response to wounding

Sensory perception: sensory perception of taste, G-protein coupled receptor protein signaling pathway, neurological process, sensory perception of chemical stimulus, sensory perception of smell

Viacej genómov pomáha vylepšiť účinnosť testov



Pozitívny výber v duplikovaných génoch



Zhrnutie

- Prirodzený výber má významnú úlohu v evolúcii
- **Purifikačný výber:**
 - Zachované regióny majú s veľkou pravdepodobnosťou nejakú funkciu
 - Pri hľadaní génov berieme do úvahy aj typické mutácie kodónov
- **Pozitívny výber:**
 - Pozitívny výber v génoch sa prejavuje veľkým pomerom nesynonymných zmien (evolúcia na proteínovej úrovni)
 - Zduplikované gény sú častejšie pod vplyvom pozitívneho výberu
 - Poľovačka pokračuje: hľadáme gény spôsobujúce charakteristické črty človeka
- **Metódy:** evolučné modely, fylogenetické HMM, test pomerov vierohodností

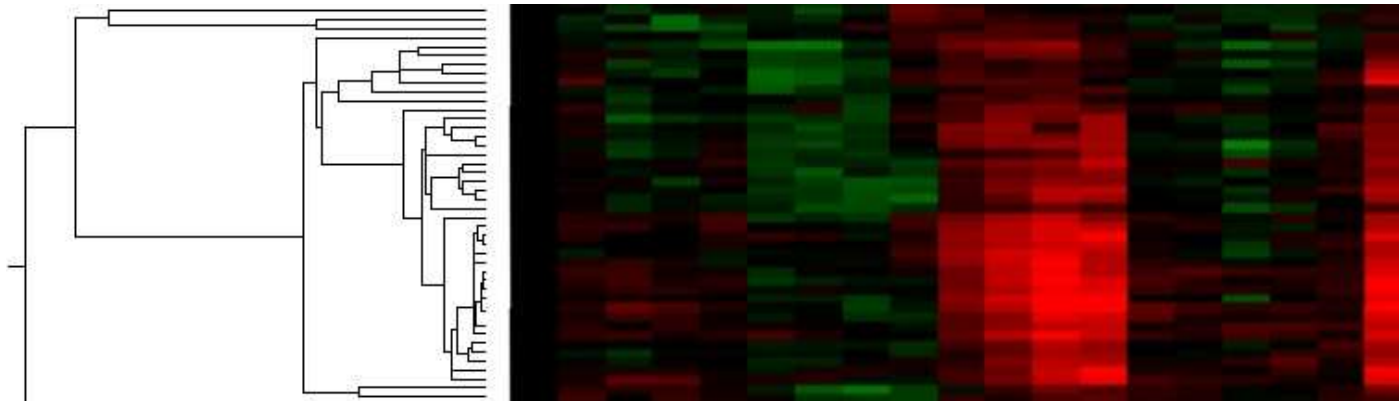
Organizačné poznámky

- DÚ 2 na stránke, odovzdať do 1.12. (na začiatku cvík, resp. prednášky)
- Stretnutia skupín 7 a 8?
- Budúci týždeň vo štvrtok sviatok, stretneme sa znovu o dva týždne

Regulácia génovej expresie

Broňa Brejová

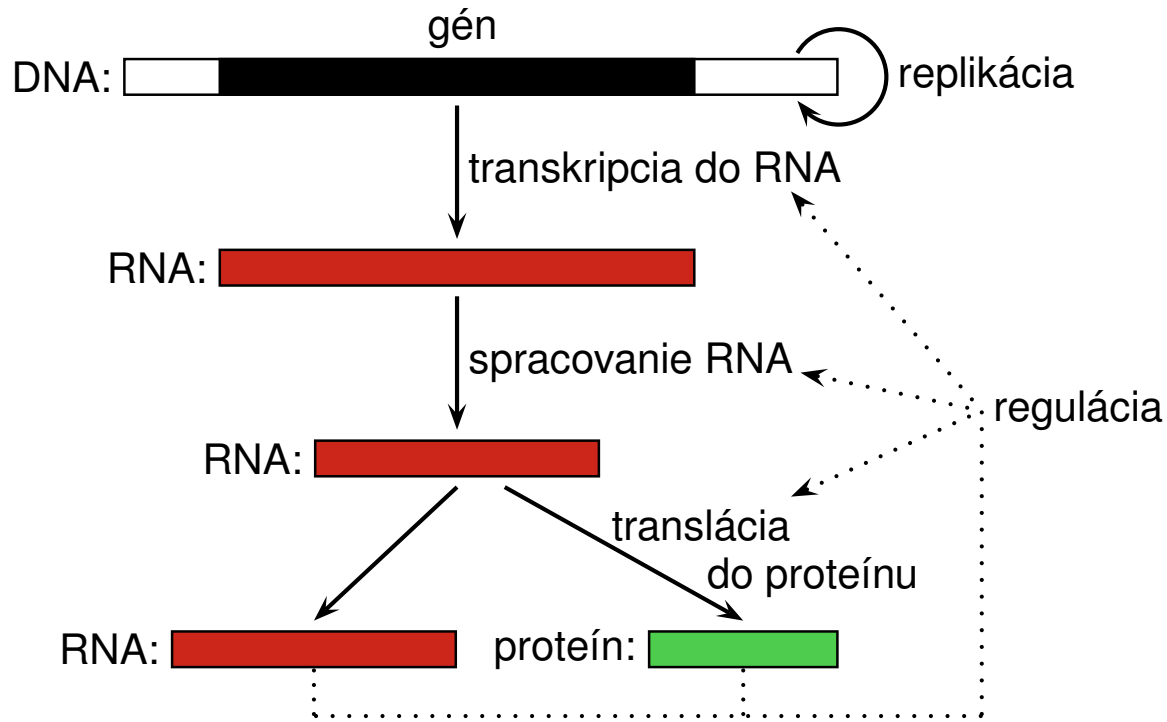
10.11.2016



Aká informácia je uložená v DNA?

Gény: Predpisy na tvorbu proteínov a funkčných RNA molekúl.

Riadenie ich expresie: kedy a koľko sa má tvoriť.

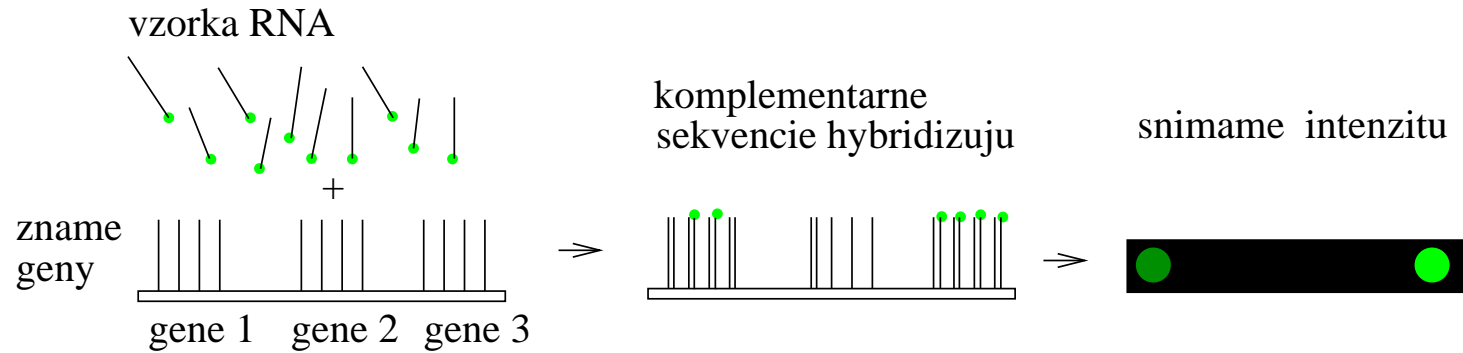


Regulácia na úrovni transkripcie, spracovania, translácie, posttranslačných modifikácií, ...

Ciele

- Zistiť, za akých podmienok je daný gén exprimovaný (súvisí s funkciou génu)
- Ktoré gény ho regulujú
- Detaily regulačného mechanizmu (väzobné miesta, zmeny v množstve expresie, . . .)

Technológia: expression array, microarray



Meranie množstva mRNA prítomnej v bunke pre **veľa génov** naraz.
Zopakujeme za rôznych podmienok.

Technológia: RNA-seq

sekvenujeme RNA extrahovanú z bunky NGS technológiami,
mapujeme na genóm, hĺbka pokrytia zodpovedá úrovni expresie



Príklad expression array dát

Pomer expresie génu v meranej a kontrolnej vzorke fg/bg

	15min	30min	1hod	2hod	4hod	...
W95909	0.72	0.1	0.57	1.08	0.66	
AA045003	1.58	1.05	1.15	1.22	0.54	
AA044605	1.1	0.97	1	0.9	0.67	
W88572	0.97	1	0.85	0.84	0.72	
AA029909	1.21	1.29	1.08	0.89	0.88	
AA059077	1.45	1.44	1.12	1.1	1.15	

...

Iyer et al 1999 The Transcriptional Program in the Response of Human Fibroblasts to Serum

Fibroblast: bunky generujúce zložky medzibunkovej hmoty

pre delenie potrebujú rastové faktory dodávané ako "fetal bovine serum"

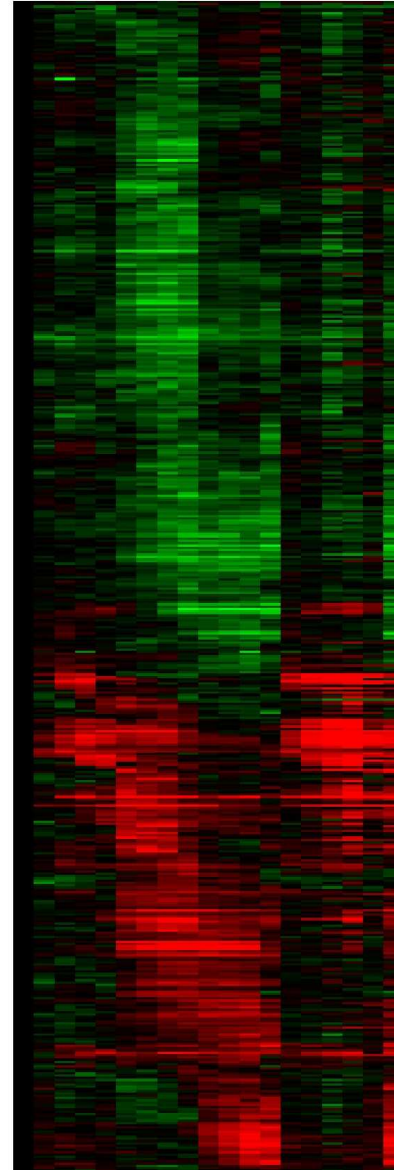
Vizualizácia

Červená: $fg > bg$

Zelená: $fg < bg$

517 génov (z 8600)

19 experimentov



Dnes: iný typ dát

- tabuľka čísel
- typické dáta v štatistike
- možno použiť všeobecné metódy štatistiky, strojového učenia

Všetky ostatné prednášky: pracujeme so sekvenciami

- zostavovanie genómov
- zarovnávanie sekvencií
- hľadanie génov
- fylogenetické stromy, populačná a komparatívna genomika
- štruktúra a funkcia proteínov a RNA

Prvá sada problémov: predspracovanie dát

- Zo scanovaných obrázkov určiť intenzitu, odhaliť zlé merania
- Agregácia dát z viacerých meraní pre jeden gén
- Použitie kontrolných meraní
- Normalizácia, aby sme mali porovnateľné výsledky z rôznych experimentov

Merania z microarray nie veľmi presné, veľa šumu, rôzne zdroje chýb

Jednoduchý výsledok:

zoznam výrazne podexprimovaných/nadexprimovaných génov

napr. $fg/bg > 2$, resp. $fg/bg < 0.5$

často na ďalšiu analýzu používame iba tieto

Zhlukovanie (clustering)

Ciel: nájsť skupiny génov s podobným profilom expresie.

Ak veľa génov v skupine má rovnakú funkciu,
ďalšie gény asi robia to isté

Meranie podobnosti profilov: napr. Pearsonov korelačný koeficient

Profil génu 1: x_1, x_2, \dots, x_n , priemer \bar{x}

Profil génu 2: y_1, y_2, \dots, y_n , priemer \bar{y}

$$C(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Číslo od -1 do 1, 1 pre lineárne korelované dáta

Vzdialenosť $d(x, y) = 1 - C(x, y)$

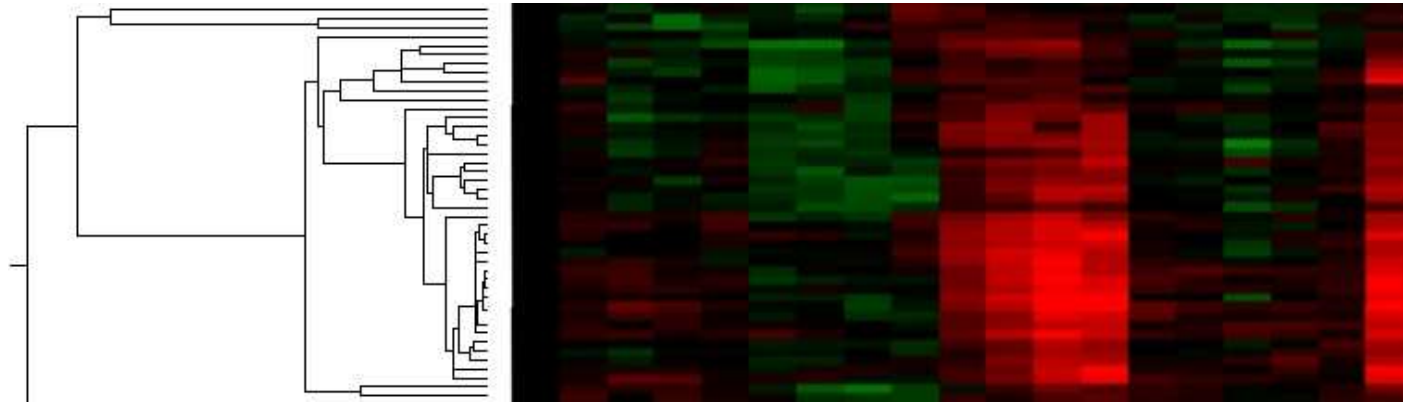
Aj iné možnosti, napr. Euklidovská vzdialenosť

Hierarchické zhlukovanie

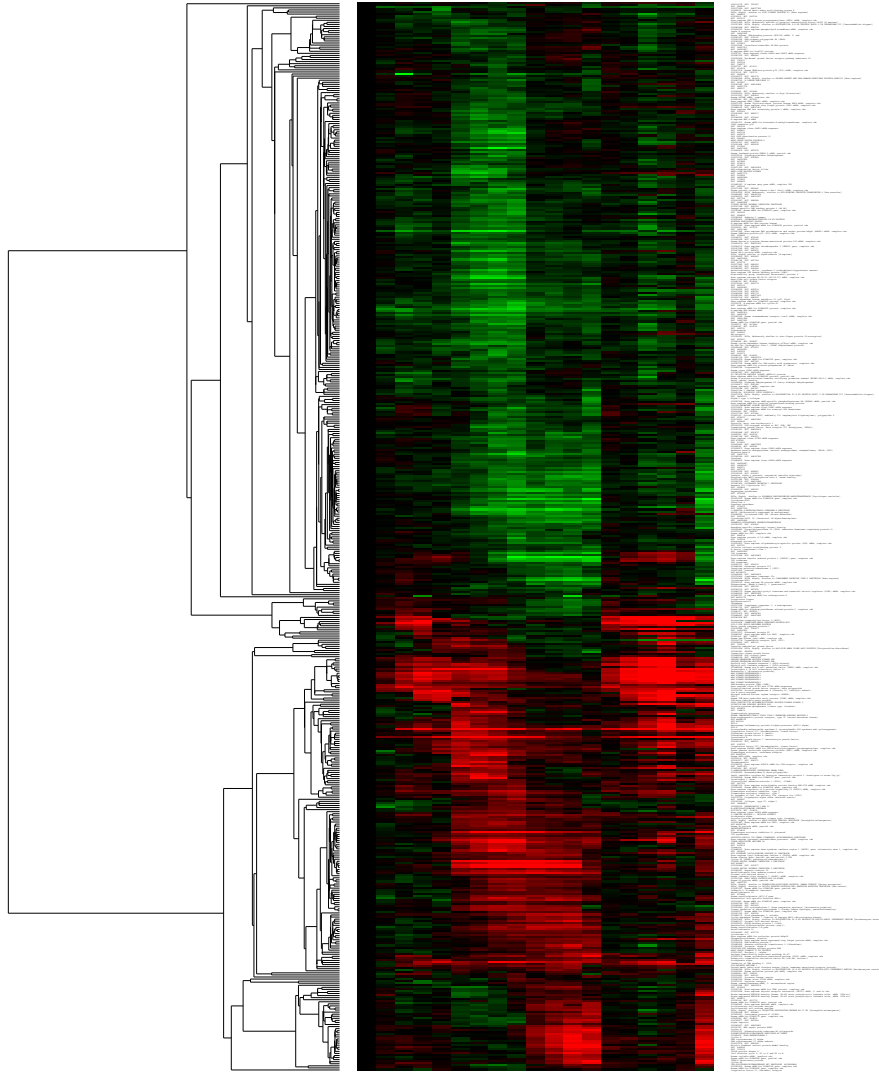
- Podobné na metódu spájania susedov vo fylogenetických stromoch
- Začneme s každým génom v samostatnej skupinke
- Nájdeme dve najbližšie skupinky a spojíme ich do jednej
- Opakujeme, kým nie sú všetky gény spolu
- Vzdialenosť skupiniek: napr. vzdialenosť najbližších génov z jednej a druhej, alebo priemer vzdialeností cez všetky páry
- Výsledkom je strom zobrazujúci postupnosť spájania

	A	B	C	D	E
gén A	0	0.6	0.1	0.3	0.7
gén B	0.6	0	0.5	0.5	0.4
gén C	0.1	0.5	0	0.6	0.6
gén D	0.3	0.5	0.6	0	0.8
gén E	0.7	0.4	0.6	0.8	0

Príklad



Zhlukovanie tiež pomáha vizualizácii dát,
podobné gény sa dostanú ku sebe

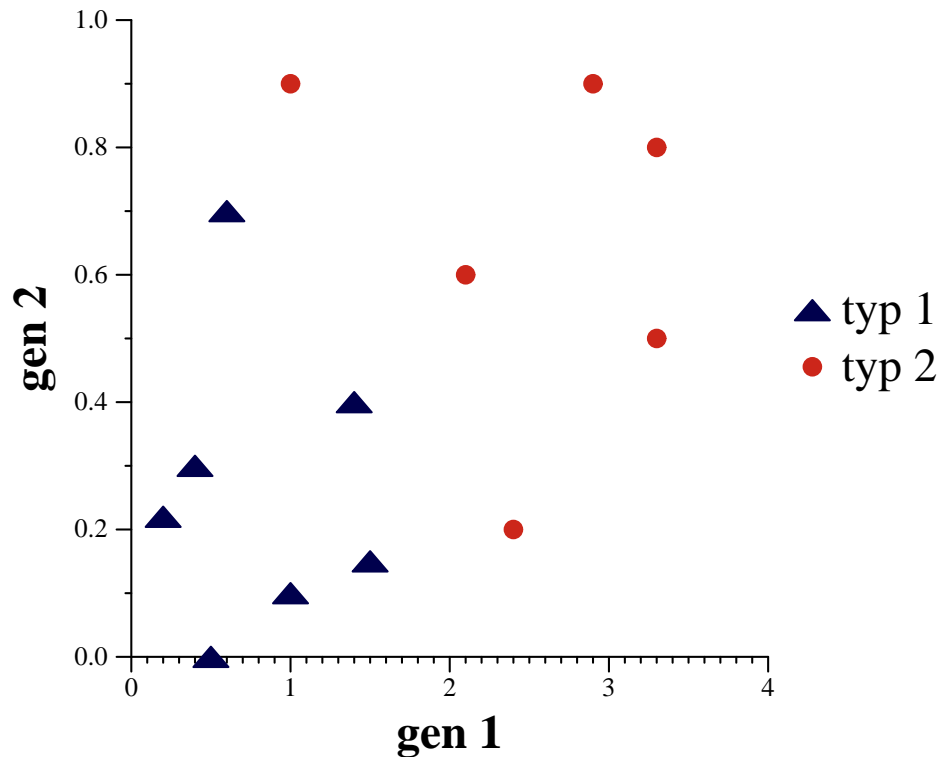


Klasifikácia

- Typický problém v strojovom učení
- Chceme odlíšiť napr. rôzne typy tumorov podľa expresie génov
- Máme nejaké príklady, kde vieme expresiu aj typ tumoru
- Chceme napr. nájsť vzorec, ktorý nám z expresie vyráta záporné číslo pre typ 1, kladné číslo pre typ 2.
- Vopred si vyberieme si typ vzorca s neznámymi parametrami (trieda hypotéz)
- Na tréningových dátach hľadáme hodnoty parametrov, pre ktoré vzorec najlepšie funguje
- Fungovanie vzorca testujeme na testovacích dátach (nepoužité na tréning)
- Hotový vzorec použijeme na dáta s neznámym typom

Jednoduchý príklad: expresia 2 génov

Trénovacie dáta so známym typom:



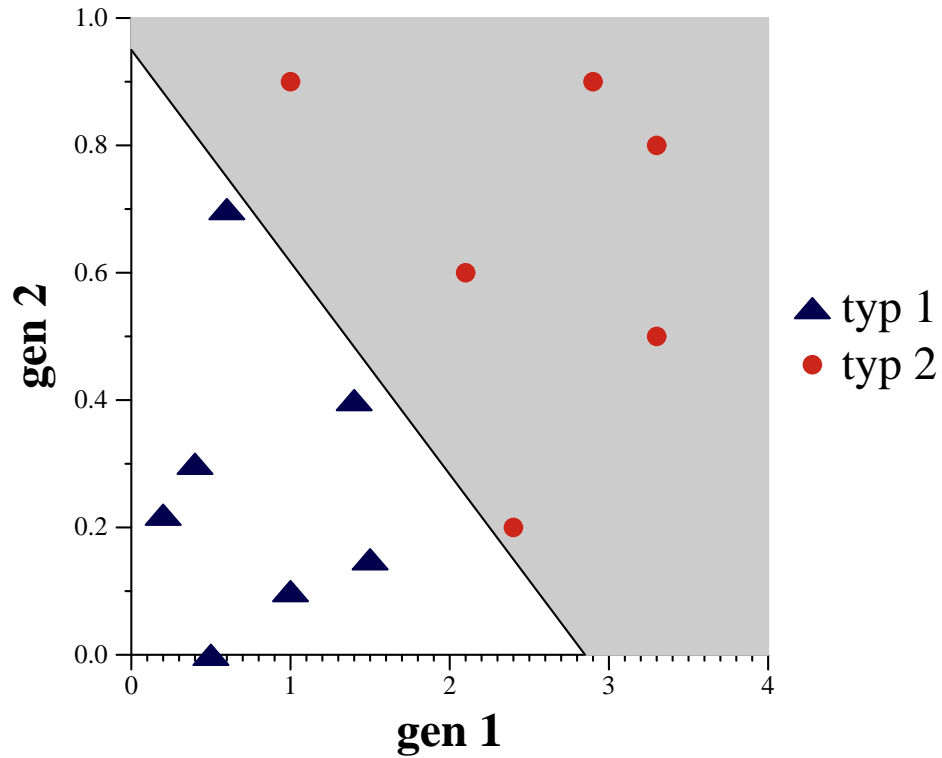
Typ vzorca: lineárne funkcie (lineárny diskriminant)

tumor typu 1 ak $ax + by + c < 0$

Hľadáme a, b, c také, aby na trénovacích dátach predpovedal dobre

Jednoduchý příklad: expresia 2 génov

Výsledný vzorec:



$$a = 1, b = 3, c = -2.85$$

tumor typu 1 ak $x + 3y - 2.85 < 0$

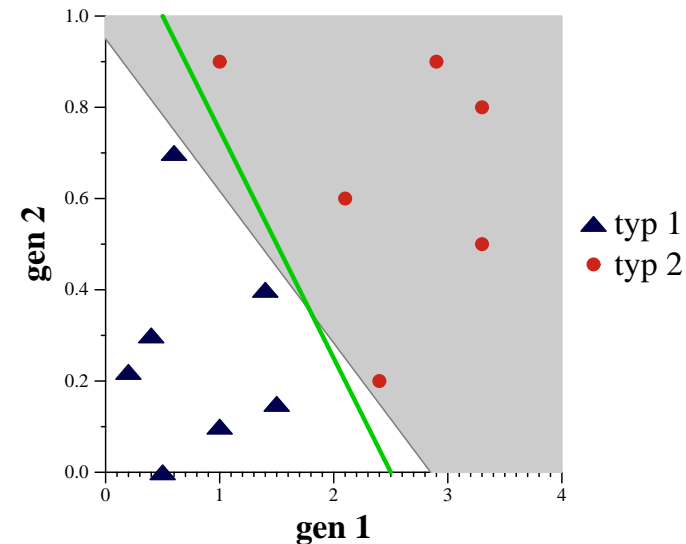
Populárne techniky na klasifikáciu

Logistic regression, logistická regresia:

lineárny diskriminátor, vracia pravdepodobnosť jednotlivých tried, dobre známa štatistická metóda.

Support vector machines

(SVM): hľadanie lineárneho diskriminátora s nulovou tréningovou chybou, ktorý je najďalej od všetkých tréningových dát.



Dá sa zovšeobecniť na nelineárne funkcie priemetom vektorov do väčšieho priestoru.

Populárne techniky na klasifikáciu

Neurónové siete:

“neuróny” poprepájané “synapsami”,
každý neurón na výstupe váhovaný priemer vstupov.

Bayesovské siete:

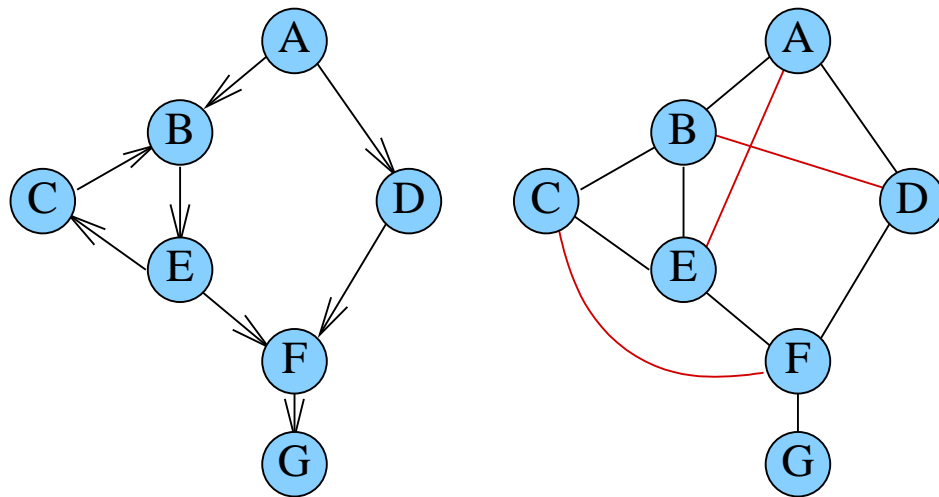
pravdepodobnostný model generujúci náhodné expresie
typ tumoru je tiež náhodná premenná, ktorej hodnotu nepoznáme
podobne ako stav v HMM

Regulačné siete z profilov expresie

Vstup: Profily expresie génov (napr. séria microarray experimentov), možno so známymi podmienkami (časové rady, delečný mutant)

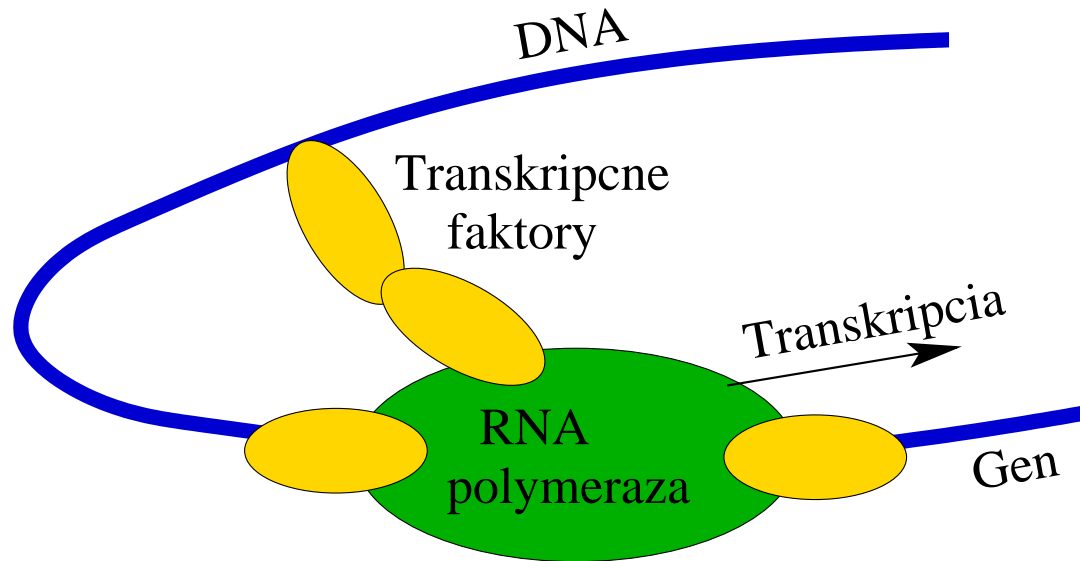
Výstup: regulačná sieť, vrcholy sú gény, orientovaná hrana $A \rightarrow B$, ak A reguluje B

Podobnosť profilov expresie nám môže dať neorientované hrany. Chceme vylúčiť hrany, ktoré vznikli tranzitivitou a správne orientovať hrany (ťažký problém)



Transkripčné faktory (TF)

Regulácia začatia transkripcie pomocou transkripčných faktorov: proteíny viažúce DNA, pomáhajú pritiahnúť RNA polymerázu

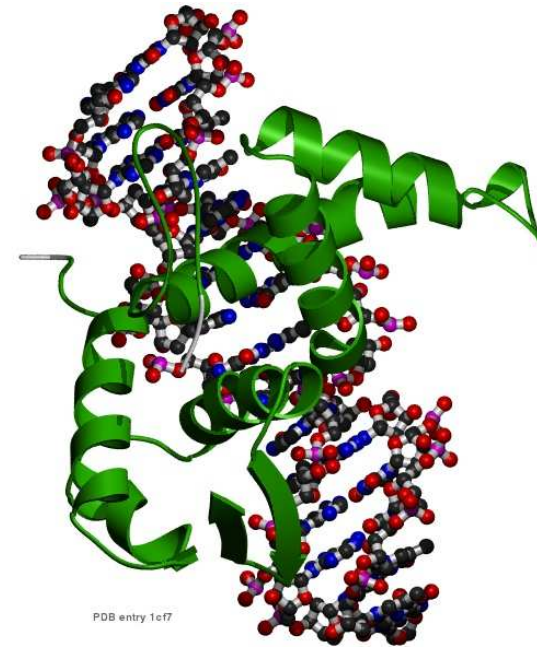
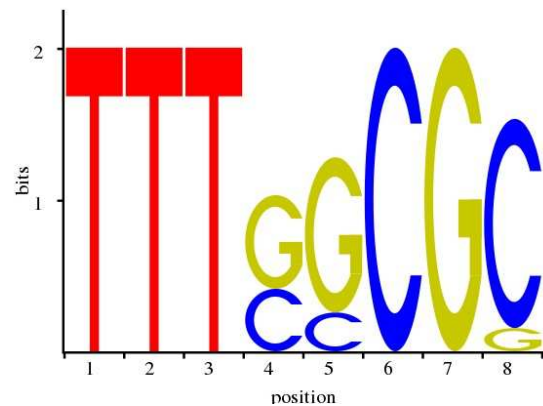


Človek má vyše 2000 TF-ov

Môžu zvyšovať alebo znižovať mieru expresie,
fungovať v skupinách

Príklad: transkripčný faktor E2F1

- Reguluje bunkový cyklus
- Viaže TTTC^{CC}GC alebo TTTC^{CG}GC, prípadne ďalšie varianty



- Sekvencie DNA, na ktoré sa viaže určitý TF chceme **reprezentovať** ako sekvenčný **motív** a hľadať **ďalšie výskyty** v genóme

Reprezentácia väzobných motívov

Reťazec s nezhodami (konsenzus):

motív je reťazec, výskyty môžu mať vopred ohraničený počet nezhôd

Príklad: motív TTTGGCGC + 1 nezhoda

TTTGGCGC, TT**A**GGCGC, TTTG**C**CGC sú výskyty motívu

TTT**C**CGC nie je výskyt

Zostavenie motívu: napr. vezmi najčastejšie písmeno na každej pozícii

A	0	0	0	0	0	0	0	0
C	0	0	0	4	2	10	0	9
G	0	0	0	6	8	0	10	1
T	10	10	10	0	0	0	0	0

Reprezentácia väzobných motívov 2

Regulárny výraz:

niektoré pozície motívu dovoľujú výber z viacej možností

[GC] znamená pozíciu, na ktorej môže byť G alebo C

N znamená hociktorú bázu

Príklad: motív TTT[CG][CG]CGC

TTTGGCGC, TTT**CC**CGC, TTTG**C**CGC sú výskyty motívu

TT**A**GGCGC nie je výskyt

Zostavenie motívu: povol' najčastejšie bázy na každej pozícii

A	0	0	0	0	0	0	0	0
C	0	0	0	4	2	10	0	9
G	0	0	0	6	8	0	10	1
T	10	10	10	0	0	0	0	0

Reprezentácia väzobných motívov 3

Position specific scoring matrix (PSSM, PWM):

skórovacia matica, skóre pre každú bázu na každej pozícii

Výskyty dosahujú skóre väčšie ako číslo T

Príklad: $T = 8$

A	-2.0	-2.0	-2.0	-2.0	-2.0	-2.0	-2.0	-2.0
C	-1.6	-1.6	-1.6	0.6	0.0	1.5	-1.6	1.4
G	-1.6	-1.6	-1.6	1.0	1.3	-1.6	1.5	-0.5
T	1.1	1.1	1.1	-2.0	-2.0	-2.0	-2.0	-2.0

TTT**CC**CGC je výskyt: $1.1+1.1+1.1+0.6+0.0+1.5+1.5+1.4=8.3$

TTTGG**C**GG je výskyt: $1.1+1.1+1.1+1.0+1.3+1.5+1.5-0.5=8.1$

TT**A**GGCGC nie je: $1.1+1.1-2.0+1.0+1.3+1.5+1.5+1.4=6.4$

Zostavenie matice z frekvencií: budúca prednáška

Hľadanie výskytov v genóme

- Hľadanie motívu v genóme: skús každú pozíciu, či je výskytom
- Väčšinou veľa falošných výskytov
- Vieme spočítať E-hodnotu: koľko výskytov očakávame v náhodnej sekvencii
- Napr. TTT[CG][CG]CGC sa vyskytuje v priemere raz za 30 000 báz
- Na zlepšenie špecificity hľadáme
 - zhluky väzobných miest,
 - miesta podporené experimentálne,
 - evolučne zachované
- Databázy motívov, napr. TRANSFAC, JASPAR

Ako nájsť väzobné miesta experimentálne?

Chromatin immunoprecipitation (ChIP)

Pomocou protilátky (antibody) na špecifický transkripčný faktor zistí, kde približne sa tento faktor viaže.

- Väzba medzi TF a DNA sa spevní formaldehydom
- DNA sa naseká na kusy
- Kusy, na ktorých je TF, sa zachytia na protilátke
- DNA sa izoluje a sekvenuje pomocou NGS (**ChIP-seq**) alebo detekuje pomocou expression array (**ChIP-chip**)

Problém: zistíme len približnú polohu väzobného miesta

Ako nájsť motívy výpočtovými metódami?

... ak nemáme niekoľko príkladov väzobného miesta

- Máme skupinu sekvencií, kde každá obsahuje väzobné miesto toho istého TF, ale väzobné preferencie TF nie sú známe
- Snažíme sa nájsť **čo najšpecifickejší** motív, ktorý sa vyskytuje vo všetkých týchto sekvenciách resp. sa vyskytuje častejšie, ako by sme očakávali.
- **Pôvodne:** zoberieme skupinu génov s podobným profilom expresie a teda možno regulovaných tým istým TF, hľadáme motív v oblastiach pred týmito génmi
- **V súčasnosti:** zoberieme oblasti detegované pomocou ChIP-seq okolo väzobných miest, nájdený motív použijeme na presnejšie určenie polohy väzby TF

Príklad: Consensus Pattern Problem (CPP)

Jednoduchá formulácia problému hľadania motívov

Vstup: dĺžka motívu L , reťazce (sekvencie) S_1, S_2, \dots, S_k

Výstup: motív (reťazec) M dĺžky L

a výskyt motívu v každom S_i (reťazec s_i dĺžky L)

také, že celkový počet nezhôd medzi M a s_i je najmenší možný

Príklad:

Vstup: CAAACAT, AGTAGC, TAACCA, TCTCCTC, $L = 4$

Výstup: motív TAAC

výskyty a nezhody AAAC 1, TAGC 1, TAAC 0, TCTC 2

celkový počet nezhôd 4

Riešenie CPP

NP-ťažký problém

- **Idea 1:** Vyskúšaj všetky možné motívy dĺžky L

Problém: Nepraktické — prečo?

- **Idea 2:** Vyskúšaj všetky možné podreťazce dĺžky L reťazcov S_1, \dots, S_k

Problém: Nemusí fungovať — prečo?

Ale dá sa dokázať, že cena riešenia bude najviac dvojnásobok optima (2-aproximačný algoritmus)

- **Ďalšie vylepšenie:** Skúšame všetky konsenzus sekvencie ℓ podreťazcov. PTAS (polynomial-time approximation scheme)

Praktickejší prístup k hľadaniu motívov

Pravdepodobnostný model generujúci sekvenciu S pomocou matice frekvencií báz v motíve W a frekvencie báz q mimo motívu

A	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
C	0.01	0.01	0.01	0.39	0.19	0.97	0.01	0.01	0.89
G	0.01	0.01	0.01	0.59	0.79	0.01	0.97	0.97	0.09
T	0.97	0.97	0.97	0.01	0.01	0.01	0.01	0.01	0.01

$$q(A) = 0.3, q(C) = 0.2, q(G) = 0.2, q(T) = 0.3$$

Pozícia motívu v S sa zvolí náhodne, každá báza sa vygeneruje z q alebo z jedného stĺpca W

Tento model definuje rozdelenie $\Pr(S | W)$.

Hľadanie motívov cez pravdepodobnostné modely

Vstup: dĺžka motívu L , sekvencie S_1, S_2, \dots, S_k , frekvencie q

Výstup: spoločný motív ako matica frekvencií W maximalizujúca vierohodnosť dát $\Pr(S_1|W) \cdot \dots \cdot \Pr(S_k|W)$

- Ťažký problém, používajú sa heuristické algoritmy
- Napríklad EM (expectation maximization)
- Lokálna optimalizácia, ktorá konverguje k lokálnemu maximu vierohodnosti
- Softvér: MEME

Schéma algoritmu EM

- **Inicializácia:**

Zvoľ si počiatočnú maticu W

(napr. zostavenú podľa jedného okna dĺžky L)

- **Iterácia:**

1. Prirad' každej pozícii j v sekvencii S_i váhu $p_{i,j}$, ktorá zodpovedá pravdepodobnosti, že na pozícii $S_i[j]$ začína výskyt motívu W .
2. Spočítaj W zo všetkých možných výskytov v S_1, \dots, S_k váhovaných podľa $p_{i,j}$

Iterácie zvyšujú vierohodnosť dát, kým nedojde ku konvergencii.

Skúšame veľa krát z rôznych počiatočných W

Zhrnutie

- Microarray alebo RNA-seq merajú úroveň expresie pre veľa génov naraz, ale v dátach veľa šumu
- Zhlukovanie (clustering) nájde podobné gény, nepotrebujeme o dátach vopred nič vedieť (unsupervised learning)
- Klasifikácia môže rozlišovať napr. choroby podľa expresie, potrebuje dáta so známou odpoveďou (supervised learning)
- Dáta o expresii pomáhajú zostaviť regulačné siete
- Väzobné motívy môžeme reprezentovať rôznym spôsobom (reťazec, regulárny výraz, skórovacia matica)
- Tieto motívy nie sú dosť špecifické, preto sa ťažko rozpoznávajú ich výskyty v genóme
- EM algoritmus na hľadanie nových motívov v sekvenciách

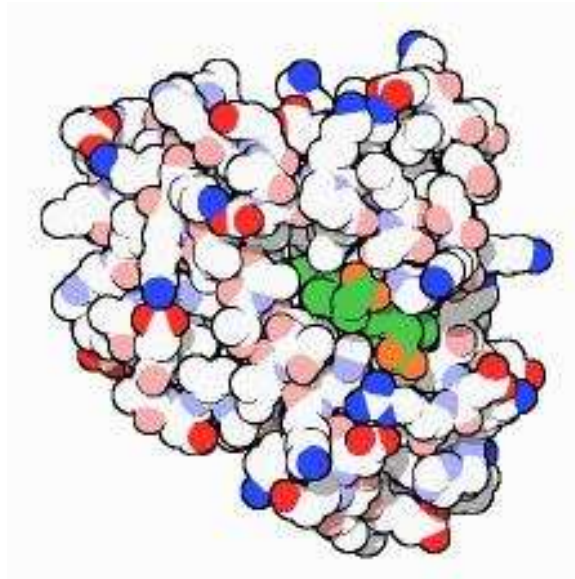
Organizačné poznámky

- DÚ2 odovzdať budúci štvrtok 1.12. (na začiatku cvík, resp. prednášky)
- DÚ3 bude zverejnená budúci týždeň, odovzdať do 15.12.
- Štvrtok 15.12. nepovinné prezentácie journal clubu
- Piatok 16.12. správy zo journal clubu
- Pondelok 12.12. návrhy projektov (ak chcete robiť projekt)
- 23.1. odovzdanie projektov (tí, čo robia projekt)
- Na budúcej prednáške bude treba základy bezkontextových gramatík, budú na dnešných cvičeniach pre biológov alebo v poznámkach z cvičení

Štruktúra a funkcia proteínov

Broňa Brejová

24.11.2016



Proteíny

Reťazce 20 rôznych aminokyselín s rôznymi chemickými vlastnosťami:

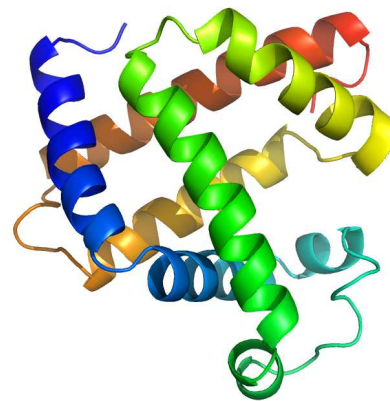
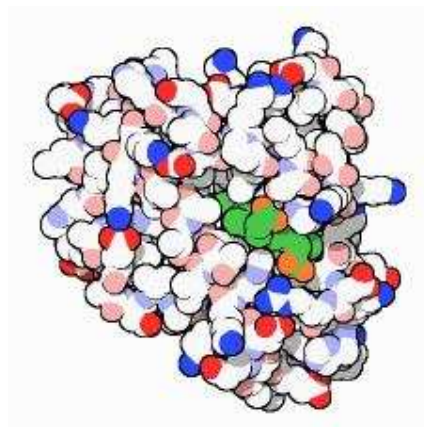
Amino Acid	Side chain	Hydrophobic	Polar	Charged	
Alanine (A)	-CH ₃	X	-	-	-
Arginine (R)	-(CH ₂) ₃ NH-C(NH)NH ₂	-	X	basic	-
Asparagine (N)	-CH ₂ CONH ₂	-	X	-	-
Aspartic acid (D)	-CH ₂ COOH	-	X	acidic	-
Cysteine (C)	-CH ₂ SH	X	-	acidic	-
Glutamic acid (E)	-CH ₂ CH ₂ COOH	-	X	acidic	-
Glutamine (Q)	-CH ₂ CH ₂ CONH ₂	-	X	-	-
Glycine (G)	-H	-	-	-	-
Histidine (H)	-CH ₂ -C ₃ H ₃ N ₂	-	X	weak basic	Aromatic
Isoleucine (I)	-CH(CH ₃)CH ₂ CH ₃	X	-	-	Aliphatic
Leucine (L)	-CH ₂ CH(CH ₃) ₂	X	-	-	Aliphatic
Lysine (K)	-(CH ₂) ₄ NH ₂	-	X	basic	-
Methionine (M)	-CH ₂ CH ₂ SCH ₃	X	-	-	-
Phenylalanine (F)	-CH ₂ C ₆ H ₅	X	-	-	Aromatic
Proline (P)	-CH ₂ CH ₂ CH ₂ -	X	-	-	-
Serine (S)	-CH ₂ OH	-	X	-	-
Threonine (T)	-CH(OH)CH ₃	-	X	weak acidic	-
Tryptophan (W)	-CH ₂ C ₈ H ₆ N	X	-	-	Aromatic
Tyrosine (Y)	-CH ₂ -C ₆ H ₄ OH	X	X	-	Aromatic
Valine (V)	-CH(CH ₃) ₂	X	-	-	Aliphatic

Štruktúra proteínov

- **Primárna štruktúra:** sekvencia aminokyselín
- **Sekundárna štruktúra:** pravidelné útvary alfa-hélix, beta-skladaný list (beta sheet)
- **Terciálna štruktúra:** presné 3D rozloženie atómov
- **Kvartérna štruktúra:** interakcia viacerých proteínov v komplexe



Myoglobín, prvý proteín so známou štruktúrou [Kendrew et al 1958]



Experimentálne určovanie štruktúry

- RTG kryštalografia (X-ray crystallography)
vyžaduje proteín v kryštalickej forme
- NMR (nuclear magnetic resonance spectroscopy)
hlavne používaná na kratšie proteíny
- Náročný a drahý proces
- Databáza štruktúr PDB
115 000 proteínových štruktúr
(UniProt má 70 miliónov sekvencií)

Bioinformatický problém: určovanie štruktúry proteínov

(protein structure prediction, protein folding)

Vstup: sekvencia proteínu

Výstup: 3D pozície atómov alebo aminokyselín

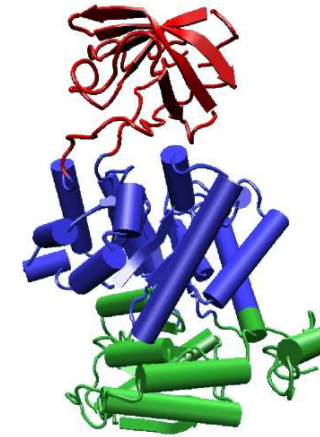
Ab initio metódy

- Nájdi štruktúru s najnižšou voľnou energiou
- Vzorce na približný výpočet energie založené na fyzike
 - sily medzi atómami v proteíne a okolitom roztoku
- Štatistické vzorce merajúce typické vzialenosti medzi aminokyselinami na známych štruktúrach
- V oboch prípadoch veľmi ťažký výpočtový problém
 - simulácia molekulárnej dynamiky
 - optimalizačné metódy, napr. simulované žihanie
- Používané na malé proteíny a zlepšenie približných štruktúr

Proteínové domény a rodiny

Doména (domain)

- Časť proteínu s nezávislou štruktúrou
- Veľa proteínov sa skladá z viacerých domén
- Domény sa tiež v proteínoch preskupujú počas evolúcie



Rodina (family)

- Skupina proteínov/domén s podobnou sekvenciou, štruktúrou, funkciou
- Ak poznáme štruktúru jedného člena rodiny, môžeme predpokladať, že ostatné majú podobnú

Proteíny ako skladačka domén

Databáza Pfam

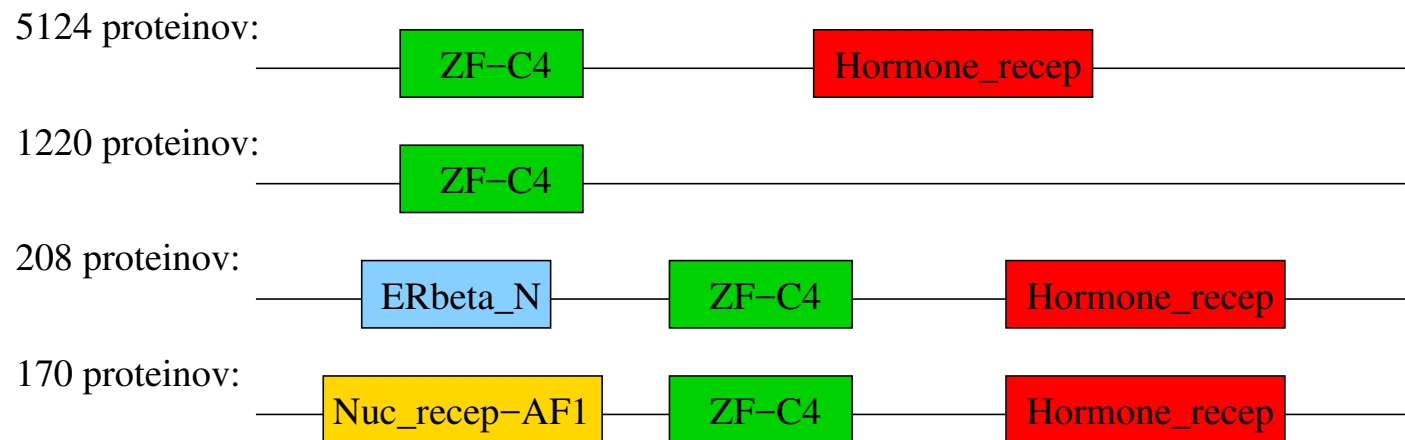
Domény v proteínoch rozdelené do rodín

80% proteínov aspoň jedna známa doména

58% proteínových sekvencií pokrývajú známe domény

Príklad:

4 z 91 architektúr obsahujúcich doménu Zinc finger, C4 type
(databáza Pfam)

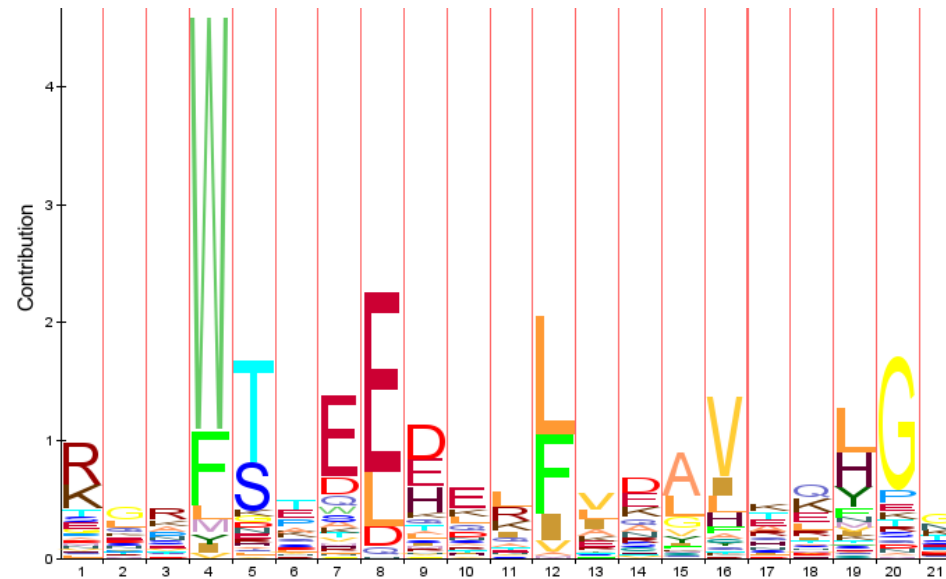


Hľadanie rodín

Cieľ: Zisti, do ktorej rodiny patrí daný proteín

- Zarovnania medzi známymi prvkami rodiny a novým proteínom nemusia nájsť vzdialených členov
- Viacnásobné zarovnanie rodiny ukáže dôležité zachované pozície
- Rodinu reprezentujeme pravdepodobnostným profilom

```
MEEWSASEANLFEEALEKYGKDF
PDEWTVEDKVLFEQAFSFHGKT.
GTKWTAENKKFENALAFYDKDT
SKNwSEDDLQLLIKAVNLFPA GT
EKPWSNQETLLLLLEAIETYGDD.
AREwTDQETLLLLLEGLEMHKDD.
KPEWSDKEILLLEAVMHYGGDD.
DDTWTAQELVLLSEGVEMYS...
KKNwSDQEMLLLLLEGIEMYE...
DENwSKEDLQKLLKGIQEF GAD.
EDDWSQAEQKAFETALQKYPKGT
EEAWTQSQQKLELALQQYPKGA
EDVwSATEQKTLEDAIKKHKSSD
AMSwTHEDEFELLKAAHKFKMG.
```



Pravdepodobnostný profil rodiny

(profile, position specific score matrix PSSM)

- V zarovnaní spočítaj $e_i(x)$: frekvencia výskytu písmena x v stĺpci i
- Dostaneme model, ktorý generuje sekvenciu x_1, x_2, \dots, x_n s pravdepodobnosťou

$$e_1(x_1) \cdot e_2(x_2) \cdots e_n(x_n)$$

- Nulová hypotéza: sekvencia bola vygenerovaná náhodne, kde písmeno x má frekvenciu $q(x)$
- Skóre: logaritmus pomeru pravdepodobností v dvoch modeloch

$$\log \frac{\prod_{i=1}^n e_i(x_i)}{\prod_{i=1}^n q(x_i)} = \sum_{i=1}^n \log \frac{e_i(x_i)}{q(x_i)} = \sum_{i=1}^n s_i(x_i)$$

Hračkářský příklad PSSM

- Uvažujme len leucín L a alanín A
- Majme zarovnanie 10 sekvencií s nasledujúcimi počtami

	1	2	3	4
A	2	6	9	1
L	8	4	1	9

- Nulová hypotéza $q(A) = 30\%$, $q(L) = 70\%$
- Sekvencia LAAL má v profile pravdepodobnosť $0.8 \cdot 0.6 \cdot 0.9 \cdot 0.9 = 0.3888$,
v nulovom modeli $0.7 \cdot 0.3 \cdot 0.3 \cdot 0.7 = 0.0441$
- Skóre $\log_2(0.3888/0.0441) = 3.14$

Hračkářsky příklad PSSM

- Majme zarovnanie 10 sekvencií s nasledujúcimi počtami

	1	2	3	4
A	2	6	9	1
L	8	4	1	9

- Nulová hypotéza $q(A) = 30\%$, $q(L) = 70\%$
- Skóre alanínu v prvom stĺpci $s_1(A) = \log_2(0.2/0.3) = -0.58$
skóre leucínu v prvom stĺpci $s_1(L) = \log_2(0.8/0.7) = 0.19$
- Dostávame tabuľku skór

	1	2	3	4
A	-0.58	1.00	1.58	-1.58
L	0.19	-0.81	-2.81	0.36

- Skóre LAAL je $0.19 + 1 + 1.58 + 0.36 = 3.13$
Skóre ALAL je $-0.58 - 0.81 + 1.58 + 0.36 = 0.55$

Pseudocounts

Ak na niektorej pozícii určitá amino kyselina nebola pozorovaná, mala by v modeli pravdepodobnosť 0

	1	2	3	4
A	2	6	9	0
L	8	4	1	10

Aby sme sa vyhli tomuto problému, pridáme ku každému políčku najskôr nejakú malú hodnotu, **pseudocount**, napr. 0,5:

	1	2	3	4
A	2.5	6.5	9.5	0.5
L	8.5	4.5	1.5	10.5

Potom postupujeme ako predtým

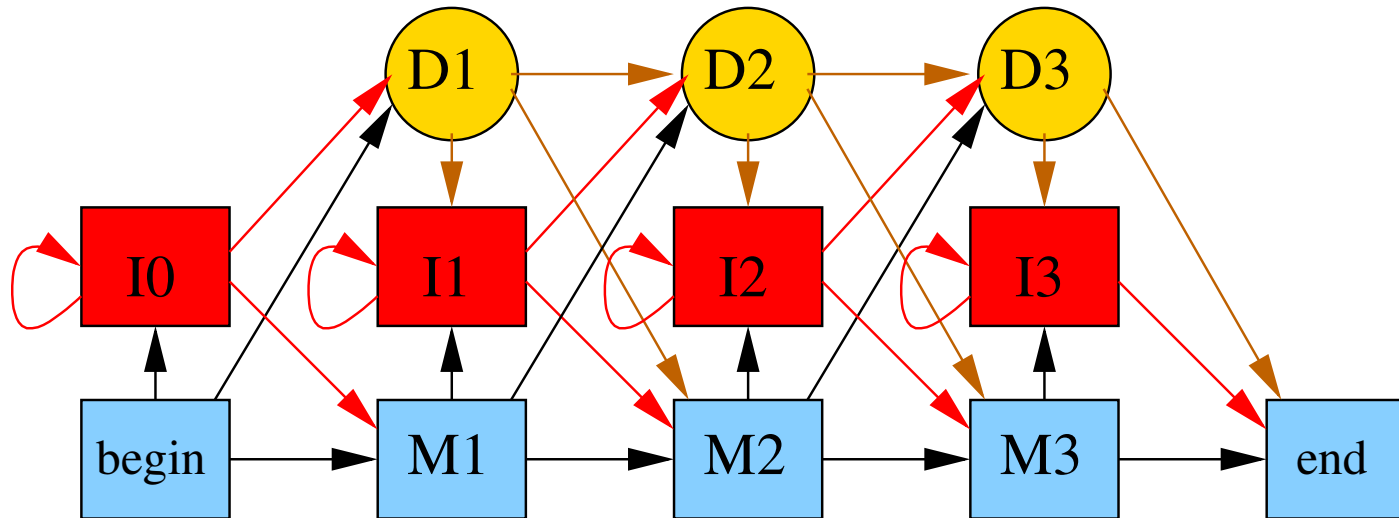
Profilové HMM

Rozšíř profil o inzercie a delécie

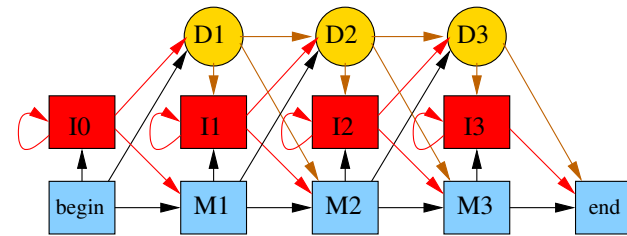
PSSM profil ako HMM:



Profilové HMM: match state, insert state, delete state



Konštrukcia profilového HMM



- Začneme s viacnásobného zarovnaní
- Stĺpcom s málo medzerami priradíme match stavy, ostatné budú v insert stavoch
- V každom stĺpci zrátame $E_i(a)$: počet výskytov a
- Pravdepodobnosť emisie $e_i(a) = \frac{E_i(a)}{\sum_b E_i(b)}$
- Pridáme “pseudocounts”, aby sme nemali nulové položky
$$e_i(a) = \frac{E_i(a)+c}{\sum_b (E_i(b)+c)}$$
- Pravdepodobnosti prechodu nastavíme podľa medzier v zarovnaní
- Veľmi podobné sekvencie môžeme použiť s menšou váhou

Použitie profilov a profilových HMM

- Pre profilové HMM používame Viterbiho algoritmus (alebo aposteriórnu pravdepodobnosť)
- PSSM profily môžeme zarovnať dynamickým programovaním s jednotným skóre pre medzery
- Rodiny domén reprezentované ako profilové HMM napr. databáza Pfam
- PSI-Blast vytvára PSSM za pochodu z podobných proteínov
- PSSM sa používajú aj na reprezentáciu motívov v DNA, napr. pre väzobné miesta transkripčných faktorov (minulá prednáška)

Protein threading

Čo ak k proteínu nenájdeme žiadnu doménu?

- Aj proteíny s pomerne odlišnou sekvenciou môžu mať podobnú štruktúru
- Môžeme skúsiť “napasovať” proteín na každú známu štruktúru
- Určitý typ zarovnania, ale pri skórovaní uvažujeme aj interakcie medzi amino kyselinami blízko v štruktúre
- Výpočtovo ťažký problém

Zhrnutie: akú štruktúru má proteín?

- Pozriem do PDB, či má známu štruktúru
- Ak nie, skúsim BLAST voči proteínom so známou štruktúrou
- Ak nič, skúsim hľadať domény so známou štruktúrou
- Ak nič, skúsim protein threading
- Pre krátke proteíny môžem skúsiť minimalizovať energiu, inak získané štruktúry doplniť/vylepšiť minimalizáciou energie

Minimalizácia energie je výpočtovo veľmi náročná

Súťaž CASP raz za dva roky

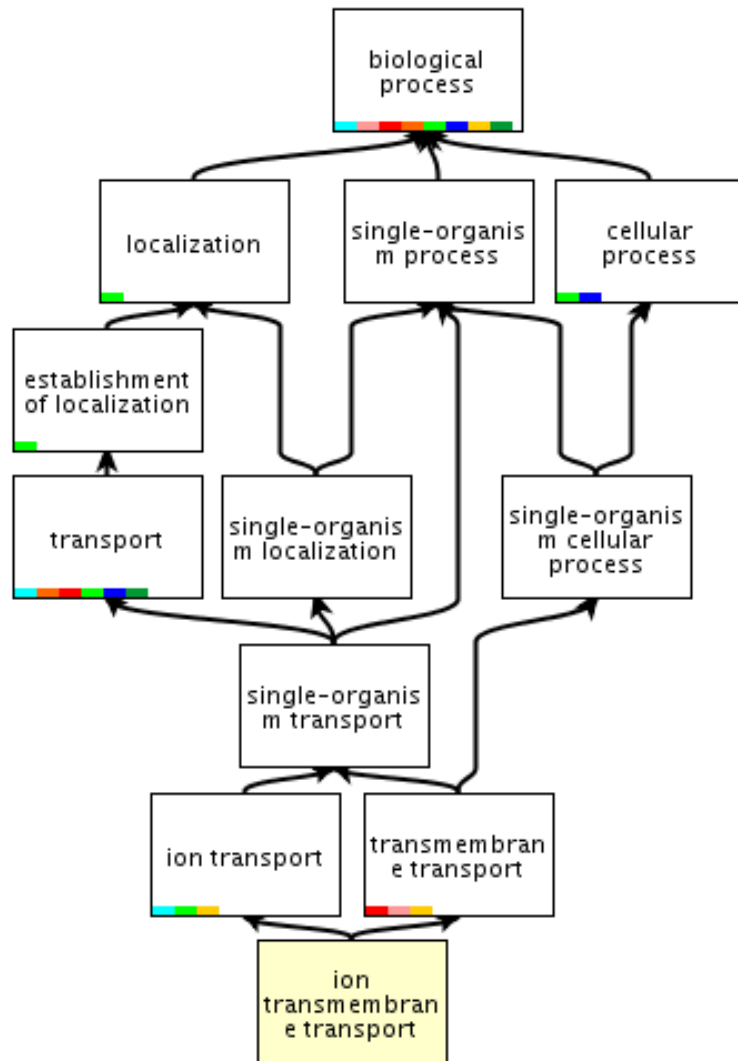
Zaujímavosti: Folding@home, Foldit

Funkcia proteínu

- Pre niektoré proteíny určená laboratórne
- Na ďalšie proteíny prenášame bioinformaticky pomocou podobnosti sekvencie, prítomnosti domén, polohy v genóme a ďalších dát
- Swissprot/Uniprot zhromažďuje údaje o funkcii proteínov
- Klasifikácia proteínov pomocou Gene ontology (GO)
Príklad pojmu v GO:
Accession: GO:0034220
Name: ion transmembrane transport
Ontology: biological_process
Definition: A process in which an ion is transported from one side of a membrane to the other by means of some agent such as a transporter or pore.
Comment: Note that this term is not intended for use in annotating lateral movement within membranes.

Gene ontology (GO)

Hierarchická štruktúra pojmov:

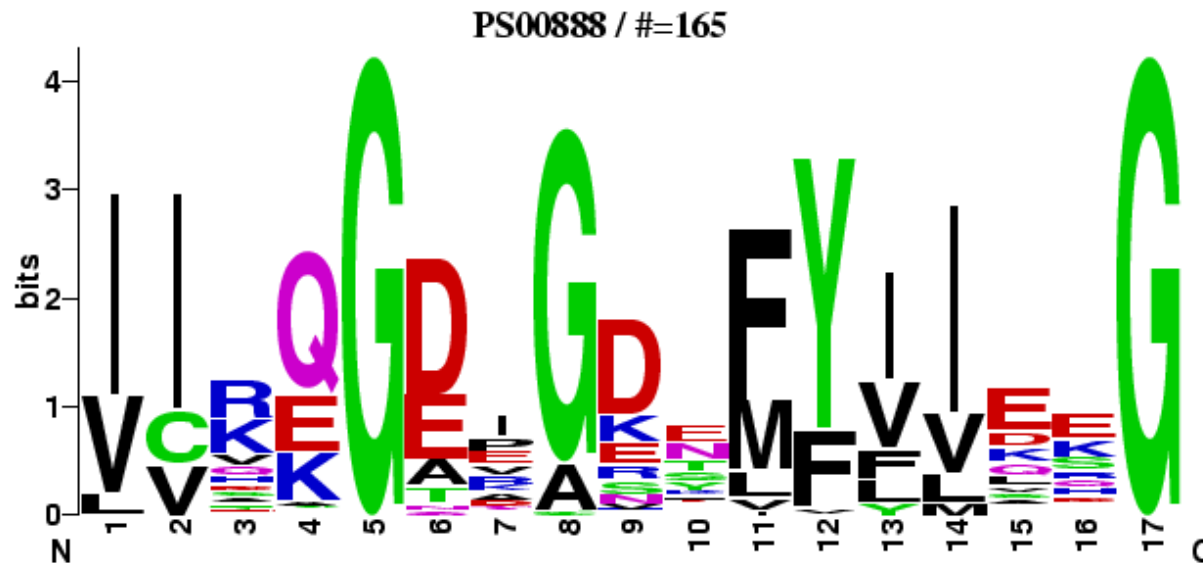


Ďalšie použitia HMM a profilov na proteíny

- Určovanie sekundárnej štruktúry
- Určovanie transmembránových proteínov a signálnych peptidov
- Určovanie funkčných motívov a posttranslačných modifikácií (databáza PROSITE)

Cyclic nucleotide-binding domain signature 1:

[LIVM] - [VIC] -x- {H} -G- [DENQTA] -x- [GAC] -{L} -x- [LIVMFY] (4) -x(2) -G



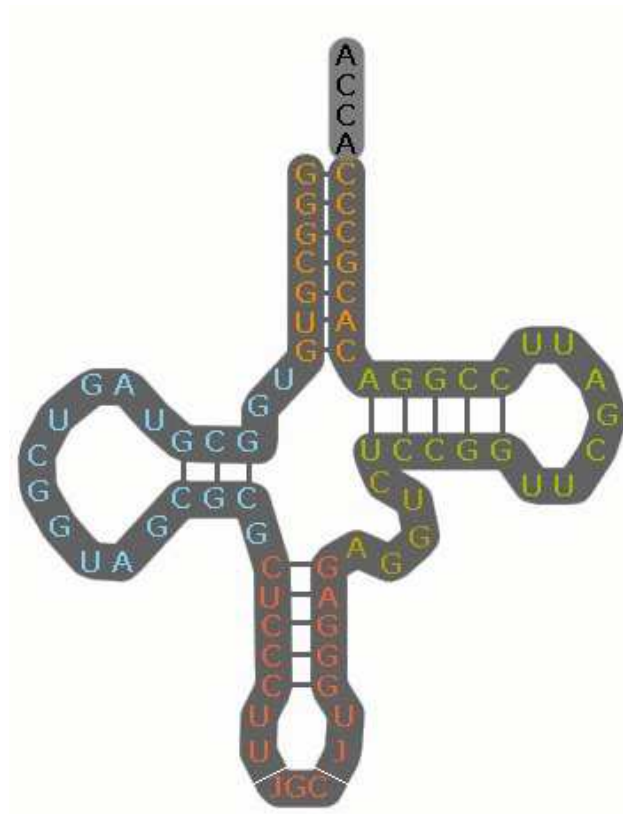
Organizačné poznámky

- DÚ 3 zverejnená zajtra ráno, odovzdať do štvrtka 15.12.
 - rovnaké zadanie pre všetkých
- Budúci štvrtok 8.12. normálny rozvrh
- Štvrtok 15.12.:
 - riadne cvičenia pre informatikov od 14:00
 - v čase prednášky nepovinné prezentácie journal clubu
 - cvičenia pre biológov nebudú
- Piatok 16.12. správy zo journal clubu
- Pondelok 12.12. návrhy projektov (ak chcete robiť projekt)
- Písomná skúška: iba jeden riadny termín
- Budúci štvrtok dohodneme:
 - či chcete prezentovať projekty (dohodnite sa v skupinách)
 - kedy bude skúška (doneste si termíny iných skúšok)

RNA

Tomáš Vinar̄

1.12.2016



Vlastnosti RNA

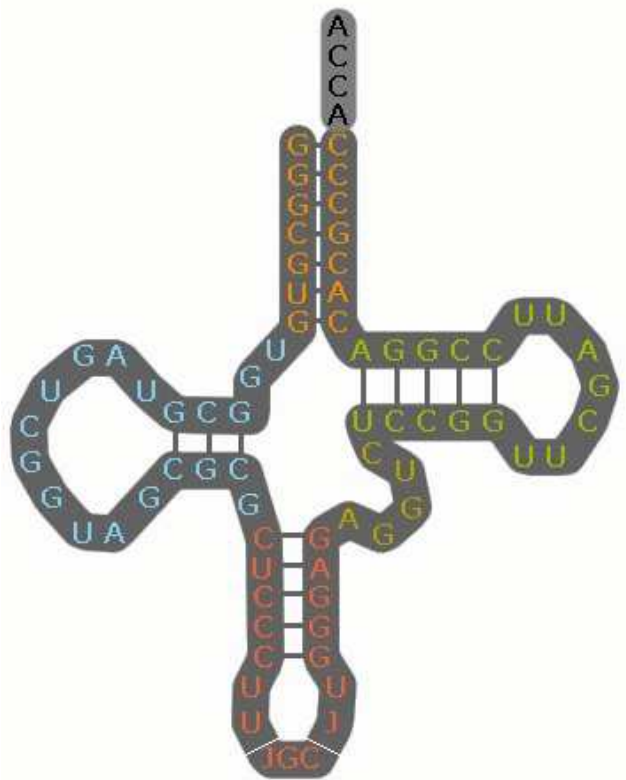
Ako sa líši od DNA?

- obsahuje ribózu namiesto deoxyribózy
- obsahuje uracil namiesto tymínu (bázy A,C,G,U)
- jednovláknové reťazce, zvyčajne kratšie
- zložitá sekundárna štruktúra: spárované komplementárne úseky
- okrem párov A-U, C-G aj nekanonické páry (napr. G-U)
- rôzne funkcie v bunke:
centrálna úloha pri expresii génov (mediátorová, transferová, ribozómová RNA),
regulácia expresie,
katalytické funkcie,
prenos genetickej informácie pre RNA vírusy

Štruktúra RNA

Príklad: transferová RNA (transfer RNA)

Sekundárna štruktúra
(secondary structure):
páry nukleotidov

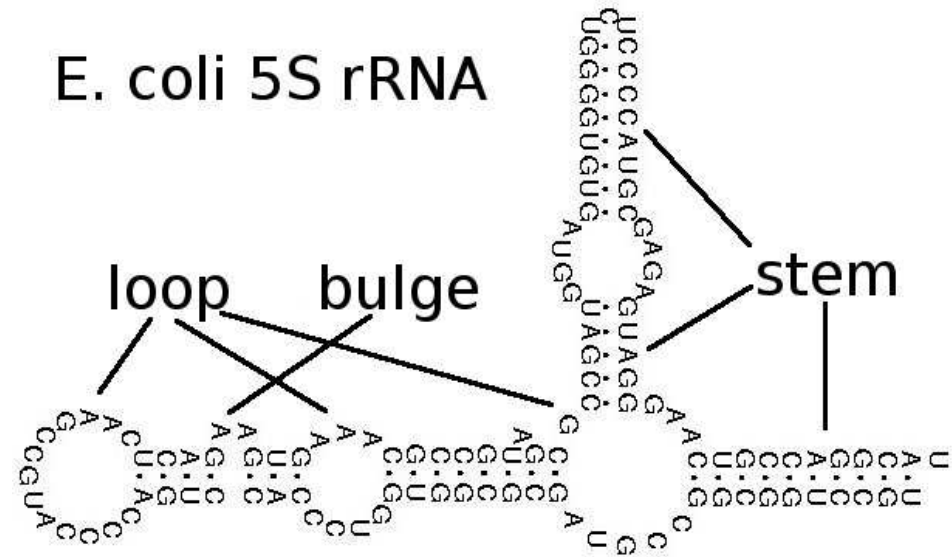


Terciárna štruktúra
(tertiary structure):
3D súradnice



Sekundárna štruktúra RNA

Prvky sekundárnej štruktúry



V tomto prípade spárované bázy tvoria **dobře uzátvorkovaný výraz**:

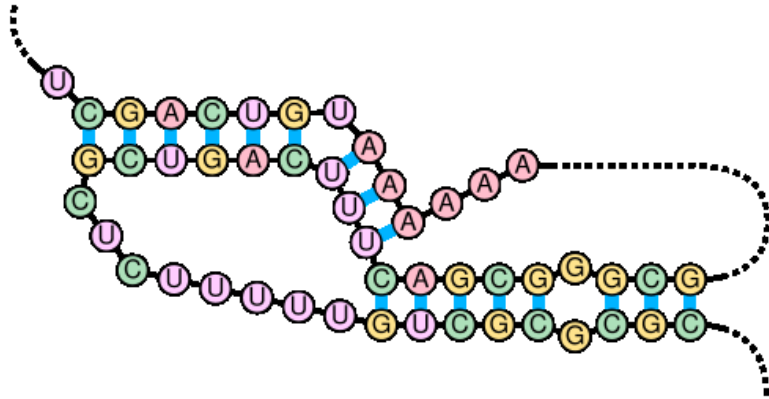
((((((((((((.....((()..)).(()..))))))))))..)).

UGCCUGGCGGCCGUAGCG...UAGCGCC...GGGAACUGCCAGGCAU

t.j. ak máme páry medzi pozíciami i a j a i' a j' a $i < i'$, tak buď $i < i' < j' < j$ alebo $i < j < i' < j'$.

Sekundárna štruktúra RNA

Pseudouzol: výnimka z dobrého uzátvorkovania



Mnohé algoritmy na prácu so sekundárnou štruktúrou ignorujú pseudouzly.

Zhruba 1.4% RNA nukleotidových párov v pseudouzloch.

Problém: určovanie štruktúry RNA

Vstup: RNA sekvencia

Cieľ: nájsť spárované bázy

Veľmi zjednodušená formulácia: nájsi dobre uzátvorkované spárovanie s najväčším počtom komplementárnych párov A-U, C-G.

Príklad: ((.((()))((.))))))

GAACACAUGUAAAUUUGUC

Možno riešiť dynamickým programovaním: [Nussinov et al., 1978]

Majme RNA X_1, \dots, X_n .

Spočítajme riešenie pre každý podreťazec X_i, X_{i+1}, \dots, X_j

$(1 \leq i \leq j \leq n)$.

Nech $A[i, j]$ je maximálny počet párov v tomto podreťazci.

Príklad: $A[1, 3] = 0$ (žiadne páry v GAA),

$A[1, 4] = 1$ (v GAAC pár G-C)

Nussinovej algoritmus

Dynamické programovanie:

Majme RNA X_1, \dots, X_n .

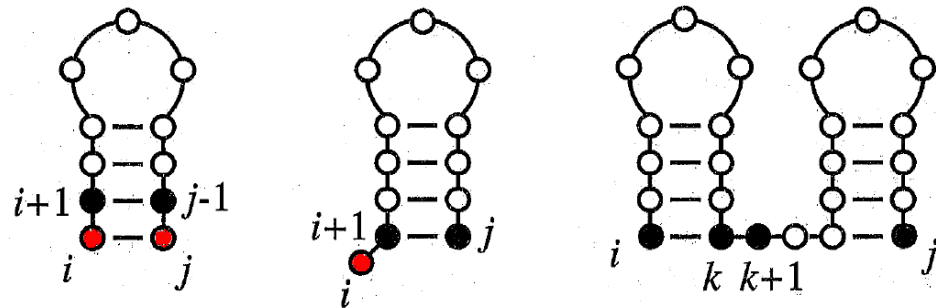
Nech $A[i, j]$ je maximálny počet párov v podreťazci X_i, X_{i+1}, \dots, X_j .

Rekurencia:

Podreťazce dĺžky 1: žiadne páry $A[i, i] = 0$

Dlhšie podreťazce: 3 prípady

- X_i a X_j sú pár: $A[i, j] = A[i + 1, j - 1] + 1$
- X_i je nespárované: $A[i, j] = A[i + 1, j]$
- X_i je pár s X_k pre $i < k < j$: $A[i, j] = A[i, k] + A[k + 1, j]$

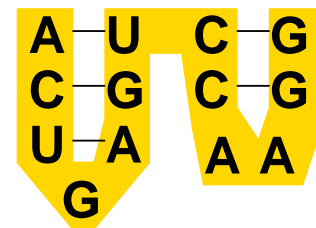


Rekurencia: $A[i, j] = \max \begin{cases} A[i + 1, j - 1] + c(X_i, X_j), \\ A[i + 1, j], \\ \max_{k=i+1 \dots j-1} \{A[i, k] + A[k + 1, j]\} \end{cases}$

	A	C	U	G	A	G	U	C	C	A	A	G	G
A	0	0	1	1	1	2	3	3	3	3	3	4	5
C		0	0	1	1	2	2	2	2	3	3	4	4
U			0	0	1	1	1	2	2	3	3	3	3
G				0	0	0	1	2	2	2	2	3	3
A					0	0	1	1	1	1	1	2	3
G						0	0	1	1	1	1	2	2
U							0	0	0	1	1	1	2
C								0	0	0	0	1	2
C									0	0	0	1	1
A										0	0	0	0
A											0	0	0
G												0	0
G													0

$$c(X_i, X_j) = \begin{cases} 1 & \text{ak } X_i - X_j \text{ môže byť pár} \\ 0 & \text{inak} \end{cases}$$

$$A[i, j] = 0 \text{ pre } i \geq j$$



Zložitosť:

$O(n^3)$ čas

$O(n^2)$ pamäť

Štruktúra s minimálnou voľnou energiou (MFE folding)

Realistickejšia formulácia problému určovania sek. štruktúry RNA.

Predpoklad: molekula v rovnovážnom stave s minimálnou Gibbsovou voľnou energiou (Gibbs free energy).

Energie pre niektoré sekvencie experimentálne zmerané.

Nearest neighbor model: sada parametrov, energie pre dvojice susedných párov v helixoch, dĺžky slučiek atď.

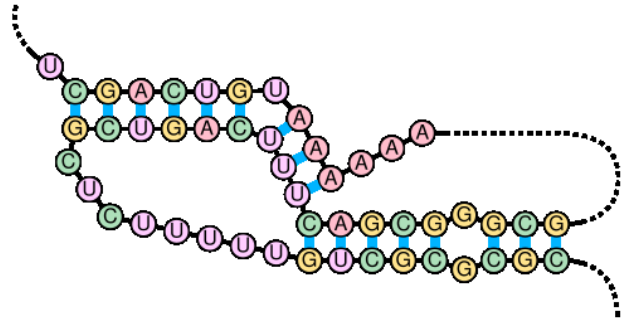
Odvodené z nameraných dát.

Príklad:

		Y:	A	C	G	U		
5'	CX	3'	-----					
3'	GY	5'	X:A		.	.	.	-2.1
			C		.	.	-3.3	.
			G		.	-2.4	.	-1.4
			U		-2.1	.	-2.1	.

Štruktúra s minimálnou energiou sa dá nájsť podobným (ale zložitejším) dyn. programovaním [Zuker and Stiegler, 1981].

Algoritmy dovoľujúce pseudouzly



Vo všeobecnosti NP-ťažký problém [Lyngso and Pedersen, 2000].

Pomalé dyn. programovanie $O(n^4)$ – $O(n^6)$ nájde niektoré typy pseudouzlov [Rivas and Eddy, 1999].

Tiež môžeme použiť heuristiky [Ren et al., 2005] (opakované vytváranie silných helixov).

Pravdepodobnostné modely na predikciu štruktúry

Chceme: model, ktorý generuje dvojice sekvencia a sek. štruktúra

Použitie: pre danú sekvenciu nájsť najpravdepodobnejšiu štruktúru

HMM nevhodné: závislosti medzi vzdialenými spárovanými bázami.

Stochastická bezkontextová gramatika, stochastic context free grammar (SCFG):

Rozšírenie bezkontextových gramatík

Pravidlám pridáme pravdepodobnosti

Stochastické bezkontextové gramatiky (SCFG)

neterminály (velké písmená) podobné na stavy v HMM,
terminály (malé písmená) reprezentují nukleotidy.

Pravidlá prepisujú neterminál na reťazec terminálov a neterminálov.
Každé pravidlo má pravdepodobnosť.

Príklad: jeden neterminál, 14 pravidiel (ϵ = prázdny reťazec)

$$S \rightarrow \begin{array}{c} 0.1 \quad 0.1 \quad 0.1 \quad 0.1 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.05 \quad 0.1 \quad 0.1 \\ \underbrace{aSu} \mid \underbrace{uSa} \mid \underbrace{cSg} \mid \underbrace{gSc} \mid \underbrace{aS} \mid \underbrace{cS} \mid \underbrace{gS} \mid \underbrace{uS} \mid \underbrace{Sa} \mid \underbrace{Sc} \mid \underbrace{Sg} \mid \underbrace{Su} \mid \underbrace{SS} \mid \underbrace{\epsilon} \end{array}$$

V každom kroku zvol' jeden (napr. najľavejší) neterminál,
prepíš ho náhodne zvoleným pravidlom:

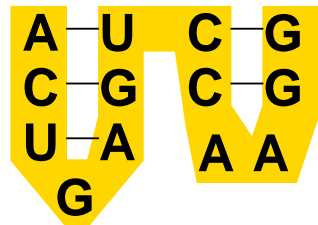
$$S \rightarrow SS \rightarrow aSuS \rightarrow acSguS \rightarrow acuSaguS \rightarrow acugSaguS \rightarrow \\ acugaguS \rightarrow acugagucSg \rightarrow acugaguccSgg \rightarrow acugaguccSagg \rightarrow \\ acugaguccaSagg \rightarrow acugaguccaagg$$

Stochastické bezkontextové gramatiky

Príklad:

$S \rightarrow aSu|uSa|cSg|gSc|aS|cS|gS|uS|Sa|Sc|Sg|Su|SS|\epsilon$

$S \rightarrow SS \rightarrow aSuS \rightarrow acSguS \rightarrow acuSaguS \rightarrow acugSaguS \rightarrow$
 $acugaguS \rightarrow acugagucSg \rightarrow acugagucgScg \rightarrow acugagucgSacg \rightarrow$
 $acugagucgaSacg \rightarrow acugagucgaacg$



Bázy vygenerované v jednom kroku sú spárované.

CYK dyn. prog. algoritmus nájde najpravdepodobnejšie odvodenie pre danú RNA v čase $O(n^3)$

Parametre možno trénovať zo známych RNA štruktúr, podobne ako pri hľadaní génov.

Gramatiky vs. minimalizácia energie

Výhody gramatík:

Parametre gramatík možno automaticky trénovať, netreba náročné experimenty.

Gramatiky sa dajú elegantne rozšíriť na modely viacerých sekvencií.

Nevýhody gramatík:

Nie je jednoduché zostaviť vhodnú gramatiku so zložitou sadou parametrov.

Nedosahujú takú presnosť ako minimalizácia energie.

Conditional log-linear models:

zovšeobecnené SCFG, tréning maximalizuje podmienenú pravdepodobnosť správnej odpovede (discriminative training).

Dosahujú lepšiu presnosť ako minimalizácia energie.

Evolúcia RNA sekvencií

Často vidíme koreláciu medzi mutáciami v spárovaných bázach. Napr. pár C-G sa zmení na G-C alebo A-U, aby sa zachovala štruktúra.

Príklad: niekoľko sekvencií z D ramena tRNA

```
(((((.....))))  
GCUCAGCC.CGGG...AGAGC  
GCCUAGCC.UGGUCA.AGGGC  
GUCUAGC...GGA...AGGAU  
GAGCAGUU.CGGU...AGCUC  
GUUCAAUC..GGU...AGAAC
```

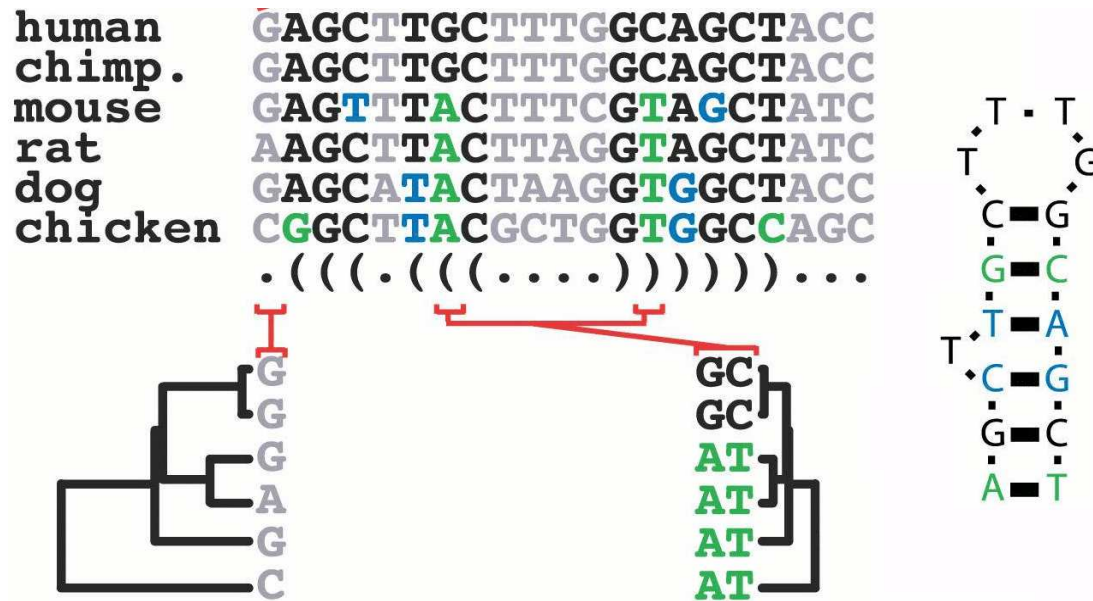
Korelácie medzi spárovanými bázami zvyšujú našu dôveru v správnosť štruktúry.

Hľadanie spoločnej štruktúry pre viacero sekvencií

- Ak sú sekvencie dostatočne podobné, môžeme ich zarovnať a potom hľadať štruktúru s veľa korelovanými párami.

Phylo-SCFG: namiesto jednotlivých báz emituje stĺpce zarovnania podľa fylogenetického stromu.

Nespárované bázy emituje bežnou substitučnou maticou, spárované bázy substitučnou maticou dvojíc (16×16).



Hľadanie spoločnej štruktúry pre viacero sekvencií

- Ak sú sekvencie dostatočne podobné, môžeme ich zarovnať a potom hľadať štruktúru s veľa korelovanými párami.
- Ak sú sekvencie málo podobné, nevieme spoľahlivo zarovnať, štruktúra však môže byť zachovaná.

Môžeme hľadať zarovnanie a štruktúru súčasne.

Presný algoritmus pomalý: $O(n^{3m})$ pre m sekvencií.

[Sankoff, 1985]

Zrýchlenie rôznymi heuristikami: predfiltrovanie, obmedzenie triedy sekundárnych štruktúr atď.

Hľadanie nových RNA génov v genóme

- Hľadaj úseky DNA so stabilnou sekundárnou štruktúrou (silnejší signál, ak máme zarovanie viacerých sekvencií).
- Výsledky treba normalizovať vzhľadom na dĺžku génu a GC%.
- Experimentálne overenie transkripcie

Problém: hľadanie známych typov RNA génov v genóme

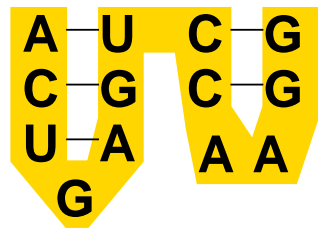
- Databáza Rfam: 2473 rodín podobných RNA génov
- Pre každú rodinu zarovnanie a pravdepodobnostný model
- Obdoba Pfamu pre rodiny proteínov
- Proteínové rodiny reprezentujeme profilmi, profilovými HMM
- Nevhodné pre RNA: závislosti medzi vzdialenými pozíciami
- Používame kovariančné modely (covariance model, CM), čo je špeciálny typ SCFG

Kovariančný model

SCFG reprezentácia RNA rodiny.

Zostavíme podľa zarovnaní + známej štruktúry.

Zjednodušený príklad



$$\begin{array}{lll}
 S \rightarrow B_1 & P_1 \rightarrow aP_2u & P_4 \rightarrow cP_5g \\
 B_1 \rightarrow P_1P_4 & P_2 \rightarrow cP_3g & P_5 \rightarrow gL_2c \\
 & P_3 \rightarrow uL_1a & L_2 \rightarrow aL_3 \\
 & L_1 \rightarrow gE_1 & L_3 \rightarrow aE_2 \\
 & E_1 \rightarrow \epsilon & E_2 \rightarrow \epsilon
 \end{array}$$

S =start, E_i =end

P_i =pár, L_i =nespárovaná báza vľavo, R_i =nespárovaná báza vpravo.

Ďalšie neterminály modelujú indely.

P_i, L_i, R_i emitujú bázy/páry s pravdepod. podľa stĺpca zarovnaní.

$$\text{Např. } P_1 \rightarrow \overbrace{aP_2u}^{0.2} \mid \overbrace{uP_2a}^{0.2} \mid \overbrace{cP_2g}^{0.4} \mid \overbrace{cP_2u}^{0.1}$$

Veľkosť gramatiky úmerná dĺžke modelovanej RNA rodiny.

Kovariančný model

Použitie:

hľadať výskyty génu v DNA (lokálne zarovnanie),
nájsť štruktúru nového génu z tej istej rodiny (globálne zarovnanie).

Dynamické programovanie: čas $O(MND^2)$,

M = počet neterminálov v gramatike, úmerný dĺžke zarovnania,

N = dĺžka DNA sekvencie,

D = max. dĺžka RNA génu v DNA (úmerná M).

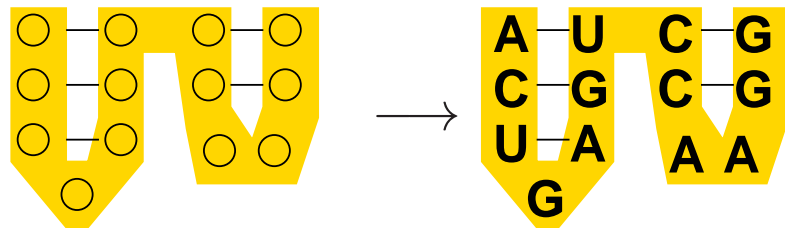
Zrýchlenie: nájsť sľubné úseky podobné na sekvencie v RNA rodine
(iba na základe podobnosti sekvencií), aplikuj CM iba na ne.

Problém: dizajn RNA

Daná RNA sekundárna štruktúra (párovanie).

Nájdí sekvenciu, pre ktorú je táto štruktúra optimálna.

Nie je známy efektívny algoritmus, heuristiky často nájdu sekvenciu pomerne rýchlo.



Použitie: skúmanie možných RNA štruktúr, vývoj liekov (ribozymes, riboswitches), RNA pre laboratórne techniky, RNA nanoštruktúry

Zhrnutie

- Určovanie sekundárnej štruktúry RNA: minimalizácia energie podľa nameraných parametrov, alebo pravdepodobnostné modely (stochastické bezkontextové gramatiky).
- Spoľahlivejšie výsledky, keď použijeme zarovnanie viacerých sekvencií, ale niekedy je ťažké správne zarovnať.
- Známe rodiny možno reprezentovať pomocou kovariančných modelov a hľadať ďalšie výskyty.
- Väčšina problémov sa dá riešiť dynamickým programovaním, ktoré je pomerne pomalé a ignoruje pseudouzly.
- Ďalšie problémy: dizajn RNA štruktúr, miRNA gény

Organizačné poznámky

- DÚ 3 do budúceho štvrtka 15.12. 14:00/15:30
- Dnes posledná prednáška
- Štvrtok 15.12.
 - cvičenia pre informatikov 14:00
 - nepovinné prezentácie journal clubu 15:30 (ktoré skupiny chcú?)
 - cvičenia pre biológov nebudú
- Piatok 16.12. 22:00 správy zo journal clubu
 - e-mailom B. Brejovej v pdf formáte, jedna za skupinu
- Písomná skúška: iba jeden riadny termín
 - kedy?
- Body z domácich úloh a journal clubu oznámime elektronicky
 - sledujte oznamy na Facebooku

Správa zo journal clubu

- Vlastnými slovami hlavné metódy a výsledky článku
- Pochopiteľná pre študentov tohto predmetu (inf aj bio)
- Netreba pokryť všetko a naopak, môžete využiť aj iné zdroje
- Skúste vložiť vlastný pohľad na tému
- Rozsah cca 1-2 strany na osobu, jeden ucelený text
- V správe vymenujte členov skupiny, ktorí sa podieľali na jej spísaní, dostanú rovnako bodov

Skúška

- Povolený ťahák na dvoch listoch A4, kalkulačka
- Na stránke ukážky jednoduchých príkladov, približne 50% bodov
 - ak niektorý neviete riešiť, spýtajte sa na Facebooku
 - v prípade pred záujmu pred skúškou konzultačné hodiny
- Zvyšné príklady budú prekvapením, v minulosti sa vyskytli:
 - Krátke príklady na pochopenie základných pojmov
 - Pre informatikov: navrhnite/modifikujte algoritmus alebo model
 - Pre biológov: usudzovanie z konkrétnych výsledkov, výber metódy pre daný problém

Polymorfizmus a populačná genetika

Broňa Brejová

8.12.2016



Populačná genetika

- Rôzne jedince toho istého druhu nemajú identický genóm
- Tieto rozdiely vplývajú na fenotyp (výzor, správanie, choroby, ...)
- Genómy viacerých jedincov môžeme sekvenovať a porovnávať s referenčnou verziou

Možné aplikácie populačnej genetiky:

- Úloha jednotlivých genetických rozdielov (škodlivé mutácie – deleterious, priaznivé mutácie – advantageous)
- História a charakter populácie (podpopulácie, migrácia, historická veľkosť populácie)

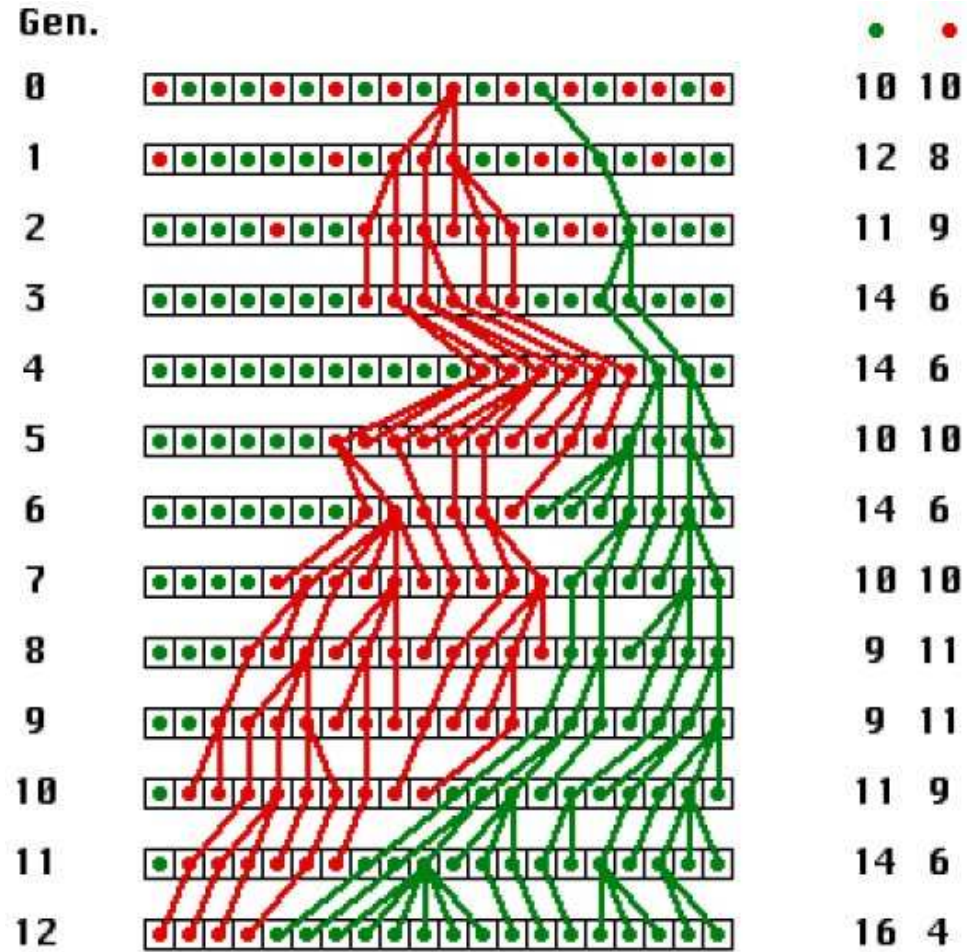
SNPy (Single Nucleotide Polymorphisms)

- Malá zmena na správnom mieste v DNA spôsobí veľké fenotypické zmeny
- SNP: jednobázová variabilita medzi jedincami ($> 1\%$ jedincov)
- Obvykle iba dve formy: **väčšinová** a **menšinová** alela

Systematické mapovanie SNPov v populácii:

- Projekt 1000 ľudských genómov: použitie NGS sekvenovania
identifikácia $> 95\%$ SNPov s aspoň 1% frekvenciou menšinovej alely

Základný model populačnej genetiky: Wright-Fisherov model



Životný cyklus SNPov vo Wright-Fisherovom modeli

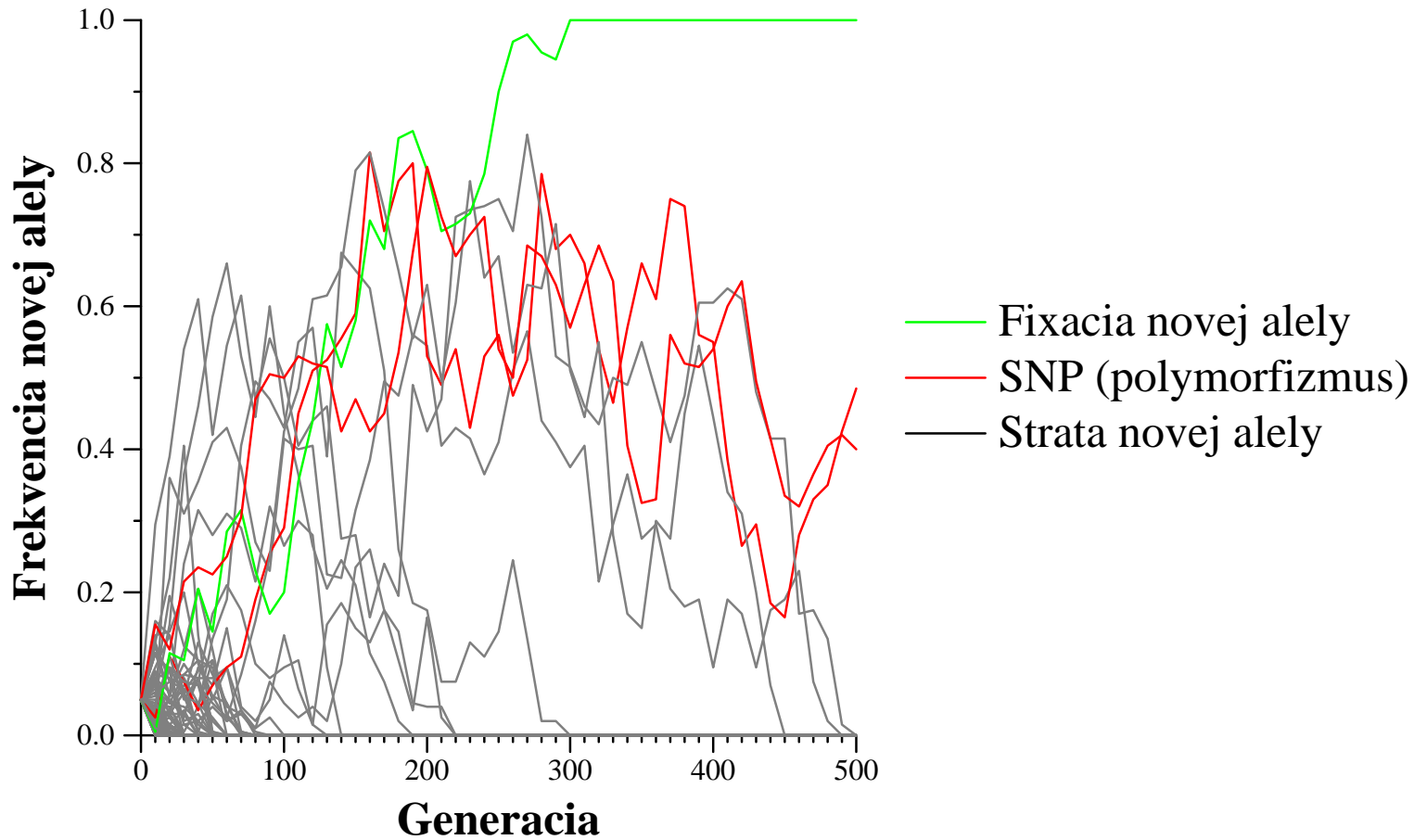
- Populácia N jedincov (stabilná veľkosť)
- Jedinec = jedna alela (A or a)
- Nová generácia vzniká “skopírovaním” náhodného rodiča (random mating), bez vplyvu prirodzeného výberu
- **Markovovský reťazec** so stavmi $\#a \in \{0, 1, \dots, N\}$

$$\Pr(\#a_t = j \mid \#a_{t-1} = i) = \left(\frac{i}{N}\right)^j \left(\frac{N-i}{N}\right)^{N-j} \binom{N}{j}$$

- Stavy 0 and N sú **pohlcajúce**

Náhodný genetický drift

$N = 200$, $X_0 = 10$, 500 generací



Zložitejšie modely populácie

- **Mutácie** zavádzajú do populácie nové alely, ktoré po čase náhodným genetickým driftom zaniknú, alebo ovládnu populáciu (fixation).
- Rýchlosť procesu je ovplyvnená efektami ako **štruktúra populácie** alebo **prirodzený výber** (selection).
- \Rightarrow Zložitejšie pravdepodobnostné modely

Mapovanie asociácií (Trait/Disease Association Mapping)

- Znaky (a choroby) vznikajú kombináciou genetických a environmentálnych vplyvov
- Cieľ: Identifikovať genetické vplyvy.
 - Ako fungujú choroby?
 - Aký je risk dedičného faktoru choroby?
 - Vývoj nových liekov, ich správne cielenie

Odbočka: diploidné genómy

- Človek je **diploidný**: má v bunke po dva chromozómy 1...22 plus pohlavné chromozómy X,X alebo X,Y
- Jeden chromozóm z páru od matky, jeden od otca
- Pre daný SNP s alelami a , A môže byť **homozygot** (aa alebo AA), alebo **heterozygot** (aA)
- Ak nejaká choroba zapríčinená alelou a , tak sa môže prejaviť iba pri homozygotoch aa , alebo aj pri heterozygotoch aA , alebo môže byť pri aa silnejšia ako pri aA
- **Haplotyp**: kombinácia aliel rôznych SNPov na tom istom chromozóme (zdedená od jedného rodiča)
Diploidný jedinec má teda dva haplotypy

Testovanie asociácie jedného SNPu

Kontingenčná tabuľka - počet haplotypov

Veľkosť psa vs. alela na pozícii chr15:44,228,468

	pôvodná alela	odvodená alela	spolu
malý pes (< 9 kg)	14	535	549
veľký pes (> 31 kg)	339	38	377
spolu	353	573	926



[Sutter a kol. 2007]

Štatisticky testujeme či sú riadky a stĺpce nezávislé (nulová hypotéza).

Ak nevyлúčime nulovú hypotézu, nepreukázali sme súvis SNPu s veľkosťou (môže ale existovať, možno treba viac dát)

Ak ju vylúčime, našli sme asociáciu, nemusí však ísť o príčinu

Testovanie nezávislosti v kontingenčnej tabuľke

	pôvodná alela	odvodená alela	spolu
malý pes	14	535	549
veľký pes	339	38	377
spolu	353	573	926

Fisherov test: (Fisher's exact test) presný výsledok z hypergeometrického rozdelenia

Chí-kvadrát (χ^2) test: obľúbený približný test, vhodný ak máme vysoké počty

Používajú sa aj zložitejšie štatistické metódy/modely (napr. diploidný genóm, príbuzenské vzťahy, ...)

Testovanie nezávislosti v kontingenčnej tabuľke χ^2 testom

	alela A	alela a	spolu
malý pes (m)	14	535	549
veľký pes (v)	339	38	377
spolu	353	573	926

V nulovej hypotéze (nezávislosť riadkov a stĺpcov) máme:

$$\Pr(A) = 353/926 = 0.381, \Pr(a) = 0.619$$

$$\Pr(m) = 549/926 = 0.593, \Pr(v) = 0.407$$

$$\Pr(A, m) = \Pr(A) \Pr(m) = 0.226$$

$$\Pr(a, m) = \Pr(a) \Pr(m) = 0.367$$

$$\Pr(A, v) = \Pr(A) \Pr(v) = 0.155$$

$$\Pr(a, v) = \Pr(a) \Pr(v) = 0.252$$

Podľa nulovej hypotézy by sme teda čakali, že 926 haplotypov bude v tabuľke rozdelených v pomeroch 0.226:0.367:0.155:0.252

Testovanie nezávislosti v kontingenčnej tabuľke χ^2 testom

Skutočná tabuľka

$O_{i,j}$ (observed):

	A	a	spolu
malý	14	535	549
veľký	339	38	377
spolu	353	573	926

Očakávané podľa nulovej hypotézy

$E_{i,j}$ (expected):

	A	a	spolu
malý	209.3	339.8	549
veľký	143.5	233.4	377
spolu	353	573	926

Spočítame veličinu $\chi^2 = \sum_{i \in \{m,v\}} \sum_{j \in \{A,a\}} \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}$

$$\chi^2 = (14 - 209.3)^2 / 209.3 + (535 - 339.8)^2 / 339.8 + (339 - 143.5)^2 / 143.5 + (38 - 233.4)^2 / 233.4 = 724.3$$

χ^2 je určitá miera rozdielnosti tabuliek O a E .

Platí, že $\chi^2 \geq 0$ a χ^2 je nula, iba ak sa tabuľky úplne zhodujú.

Testovanie nezávislosti v kontingenčnej tabuľke χ^2 testom

$O_{i,j}$ (observed):

	A	a	spolu
malý	14	535	549
veľký	339	38	377
spolu	353	573	926

$E_{i,j}$ (expected):

	A	a	spolu
malý	209.3	339.8	549
veľký	143.5	233.4	377
spolu	353	573	926

Spočítame veličinu $\chi^2 = \sum_{i \in \{m,v\}} \sum_{j \in \{A,a\}} \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}} = 724.3$

Ak platí nulová hypotéza, χ^2 je približne z rozdelenia $\chi^2(1)$,
t.j. **chí kvadrát s jedným stupňom voľnosti**.

1 stupeň: ak poznáme E a 1 políčko z O , zvyšok O vieme dopočítať.

Šanca, že pri nulovej hypotéze nám náhodne vyjde $\chi^2 \geq 724.3$ je $1.6 \cdot 10^{-159}$ (P-hodnota)

Na **odmietnutie nulovej hypotézy** často používame
prah $P < 0.05$, t.j. $\chi^2 > 3.841$

Závislosť medzi dvoma rôznymi SNPmi

Uvažujme SNP s alelami p/P a ďalší SNP s alelami q/Q .

Nameriame počty haplotypov pq, PQ, pQ, Pq

Príklad: 2000 haplotypov (1000 jedincov)

	Q	q	
P	474	611	$\chi^2 = 184.78$, P-hodnota $4.4 \cdot 10^{-42}$
p	142	773	

Stĺpce a riadky teda nie sú nezávislé, medzi SNPmi je závislosť

Príklad 2: Podobné pomery počtov, ale iba 30 haloptypov:

	Q	q	
P	7	9	$\chi^2 = 3.0867$, P-hodnota 0.07893
p	2	12	

Nulovú hypotézu nevyhlúčime pre prah $P < 0.05$ ($\chi^2 > 3.841$)

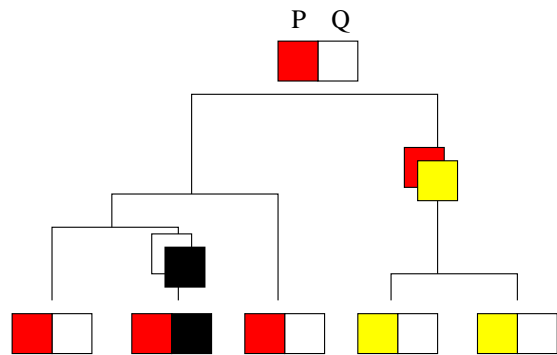
Ale pozor, pre takéto malé hodnoty χ^2 **nepresný**

Ako vzniká závislosť medzi dvoma rôznymi SNPmi

Na rozdielnych chromozómoch:

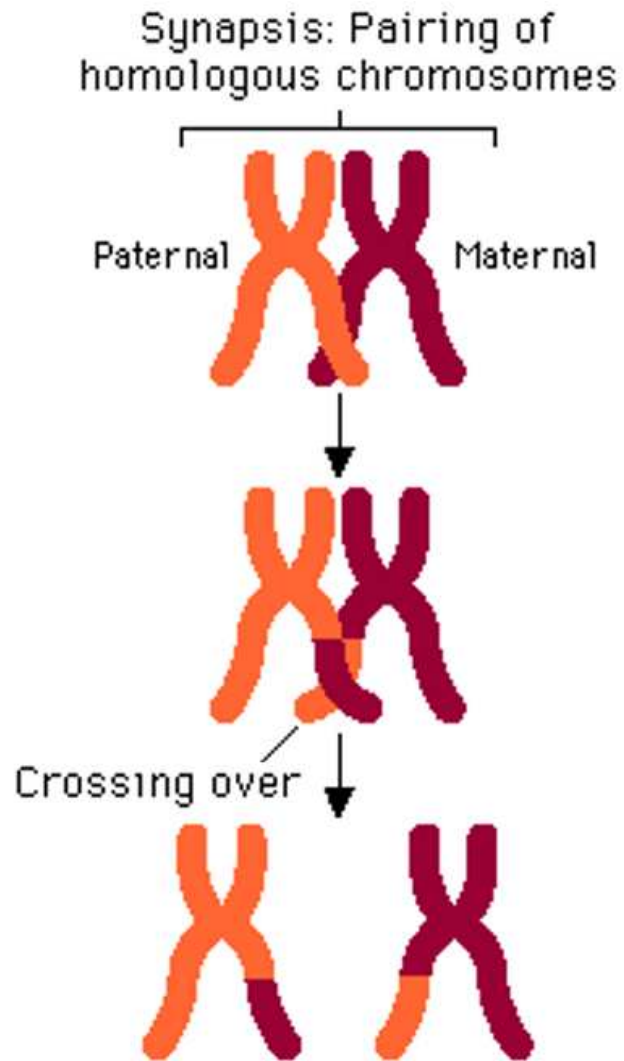
- Pravdepodobnosti výskytu jednotlivých alel sú nezávislé
- $\Pr(pq) = \Pr(p) \Pr(q)$, $\Pr(PQ) = \Pr(P) \Pr(Q)$, atď
- **väzbová rovnováha, linkage equilibrium (LE)**

Blízko seba na tom istom chromozóme:



- Málokedy mutácia na to istom mieste 2x, zriedkavá rekombinácia
- Kombinácie nie sú úplne náhodné
- Korelácie medzi SNPmi
⇒ **väzbová nerovnováha, linkage disequilibrium (LD)**

Rekombinácia



Cca 1-3 **rekombinácie** v 1 ľudskom chromozóme počas meiózy (tvorba pohlavných buniek)

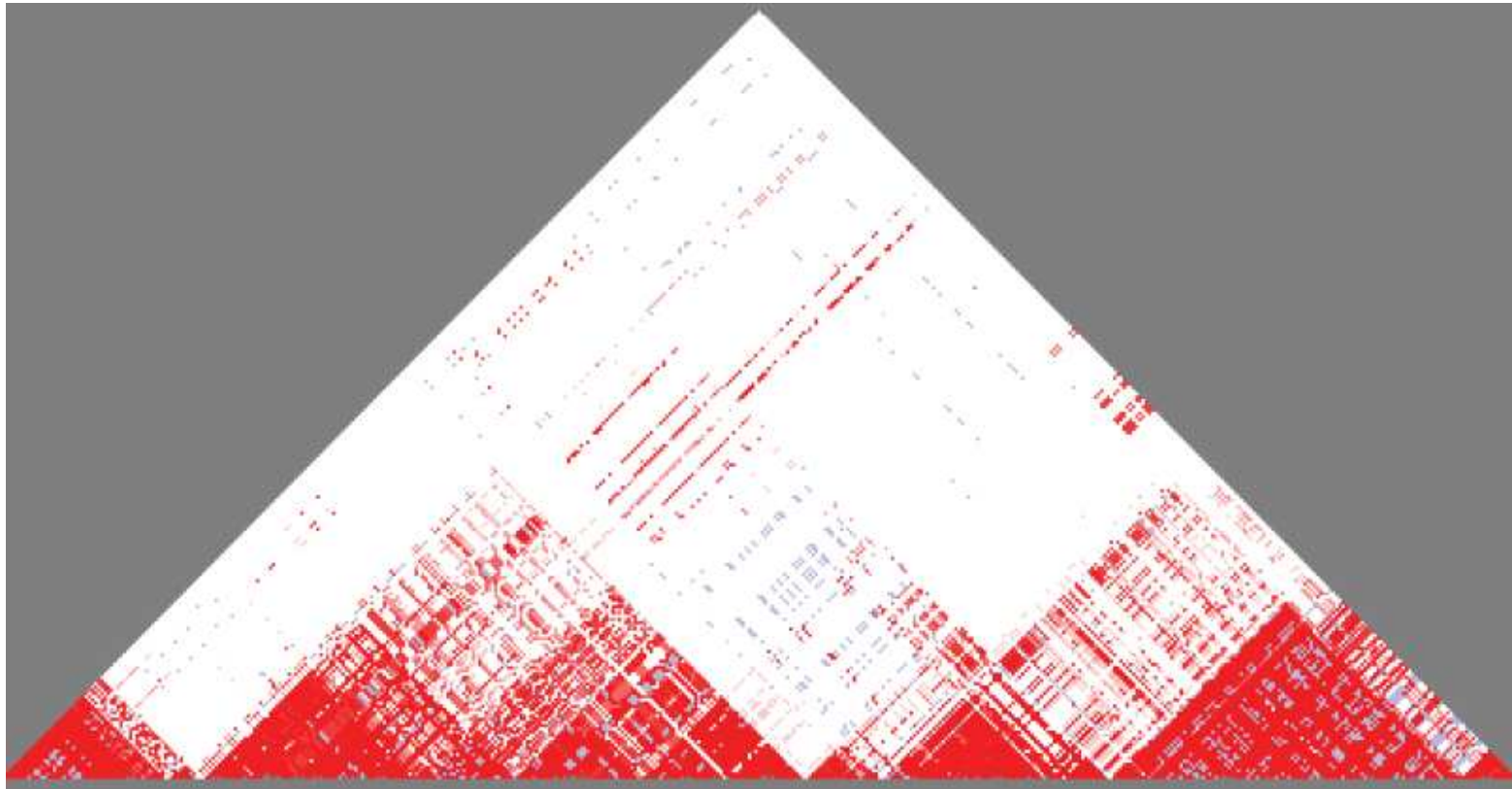
Rekombinácia znižuje LD

Ak predpokladáme rovnomernú rekombináciu:

- Čím vzdialenejšie SNPy, tým nižšie LD
- Čím staršie SNPy, tým nižšie LD
- Ďalšie aspekty: štruktúra populácie, prirodzený výber, rekombinačné hotspoty

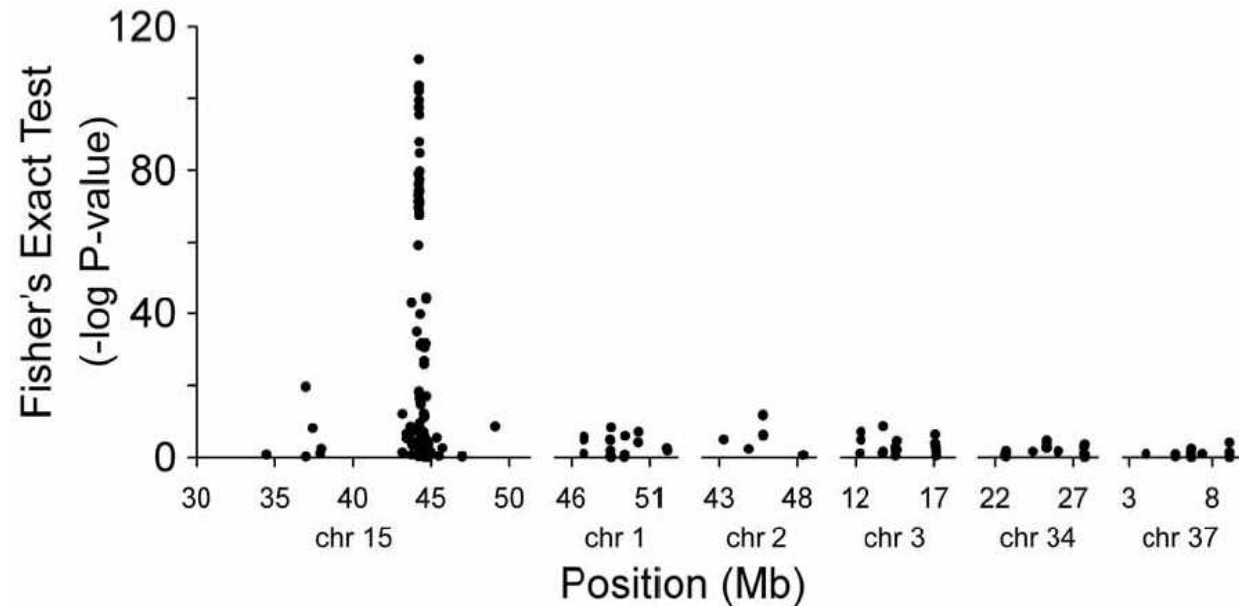
Linkage disequilibrium v ľudskom genóme

[The International HapMap Consortium, 2005]



Encode región ENm014 (500kB, chr 7), 90 ľudí Utah

Späť k psom: Hľadanie asociácií v celom genóme (Whole-Genome Association Scan, WGAS)



- V prípade štúdie veľkosti psov: WGAS identifikoval 84 kB región
- Pozíciu ďalej treba spresniť ďalšími experimentami
- **Malé LD bloky** \Rightarrow potreba veľkého rozlíšenia SNPov
- **Veľké LD bloky** \Rightarrow príliš veľké výsledné regióny

Populačno-genetické analýzy na základe pravdepodobnostných modelov

Typické parametre pravdepodobnostného modelu:

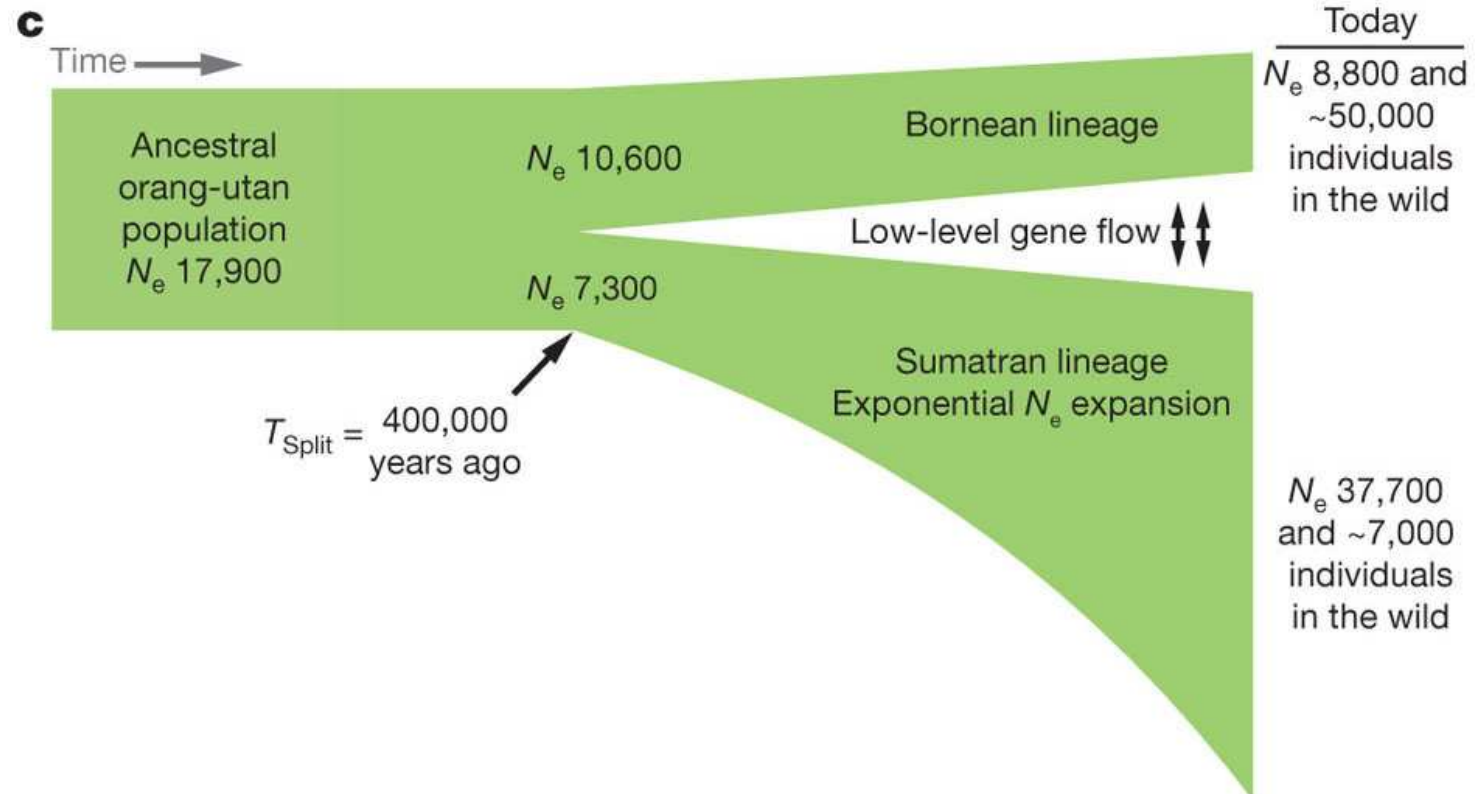
- efektívna veľkosť populácie
- rýchlosti rekombinácie a mutácie

Parametre ovplyvňujú pozorované dáta:

- Frekvencie SNPov (frekvencia menšinovej alely)
- Heterozygocita u diploidných jedincov
- Počet a veľkosť LD blokov

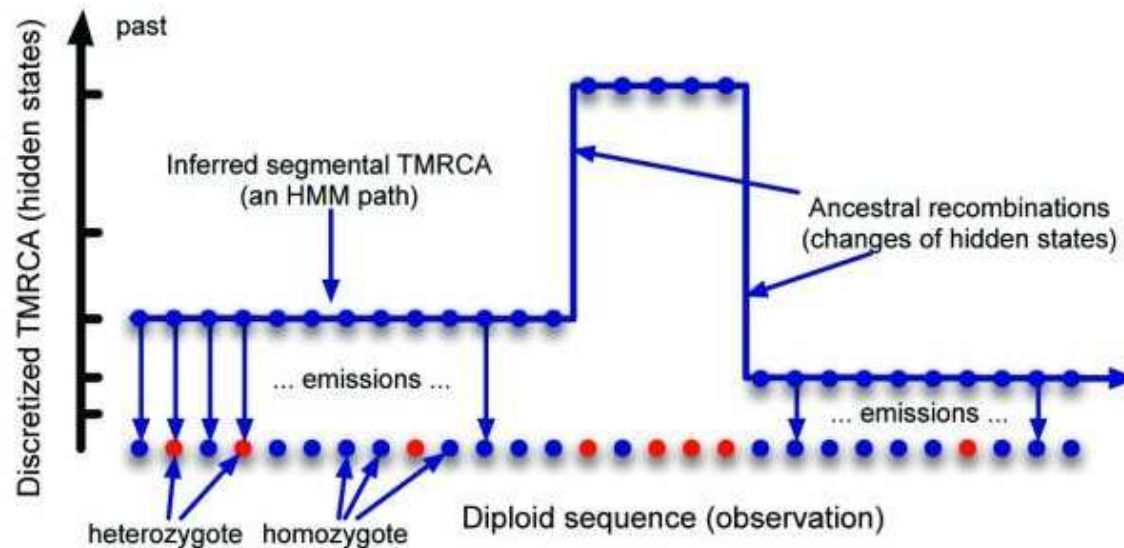
Štandardný prístup: Snažíme sa nájsť parametre modelu, ktoré najlepšie vysvetľujú pozorované dáta u osekvenovaných jedincov.

Príklad: Populačná história orangutánov



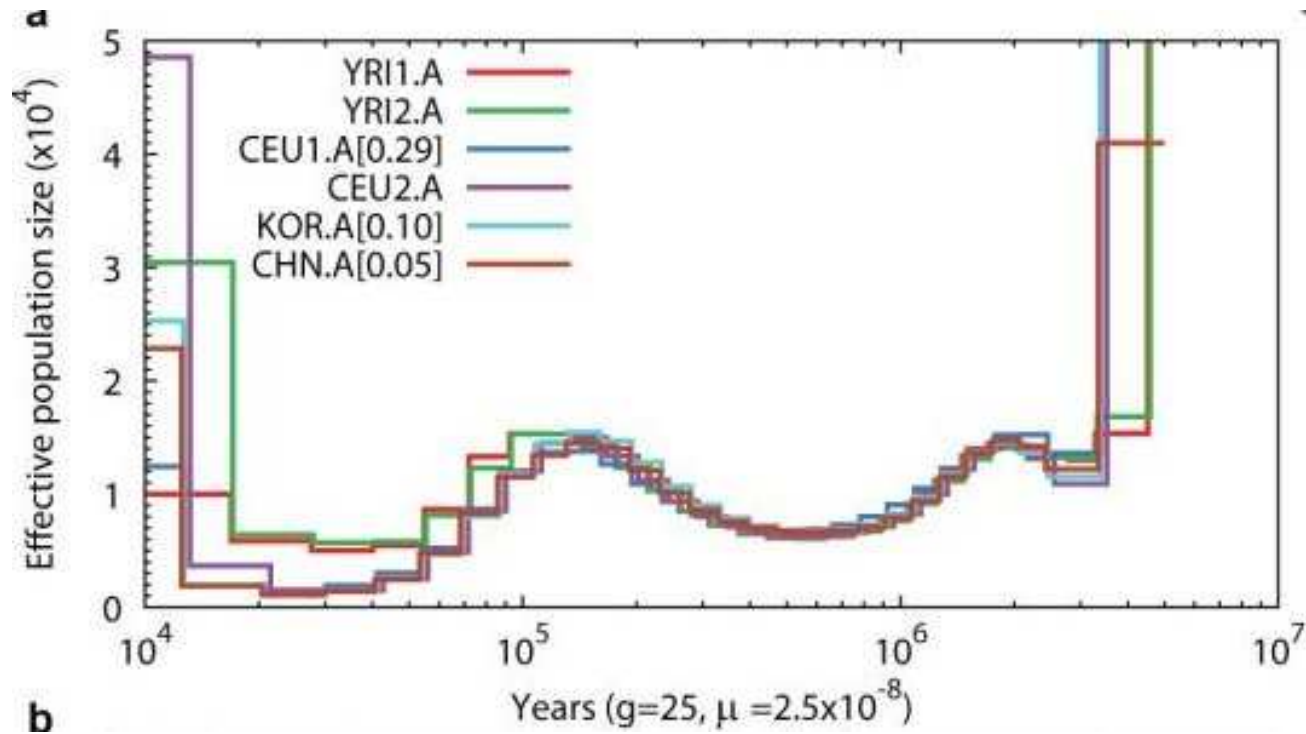
Príklad: História ľudskej populácie z genómu jedinca

- **Parametre modelu:** história vývoja efektívnej veľkosti ľudskej populácie v čase
- **Pozorované štatistiky:**
 - rozdelenie veľkostí rekombinačných blokov
 - rozdelenie časov ku najbližšiemu spoločnému predkovi (TMRCA)



Príklad: História ľudskej populácie z genómu jedinca

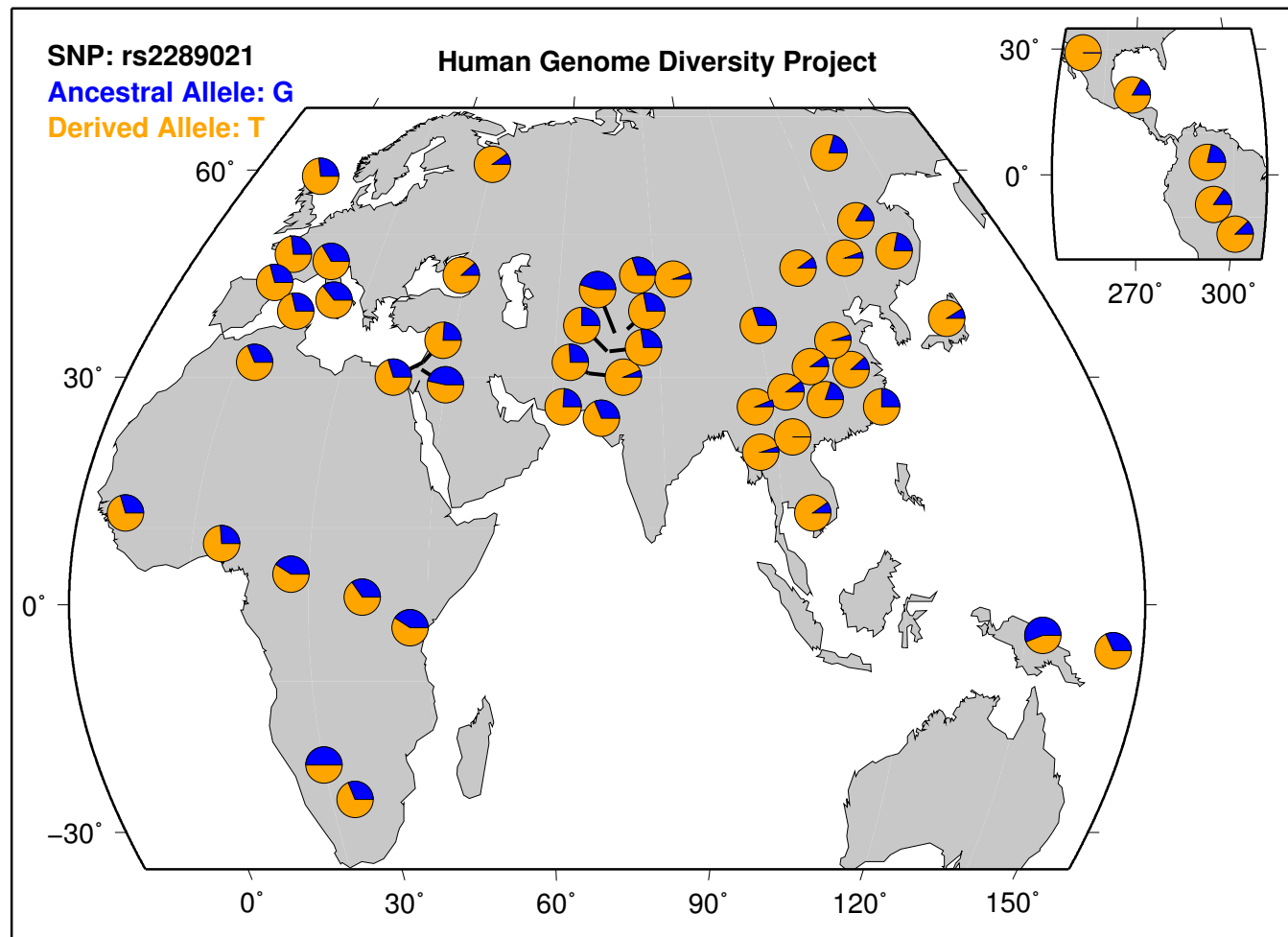
Úloha: Nájdi históriu vývoja efektívnej veľkosti ľudskej populácie, ktorá najlepšie vysvetľuje pozorované štatistiky



Štruktúra populácie

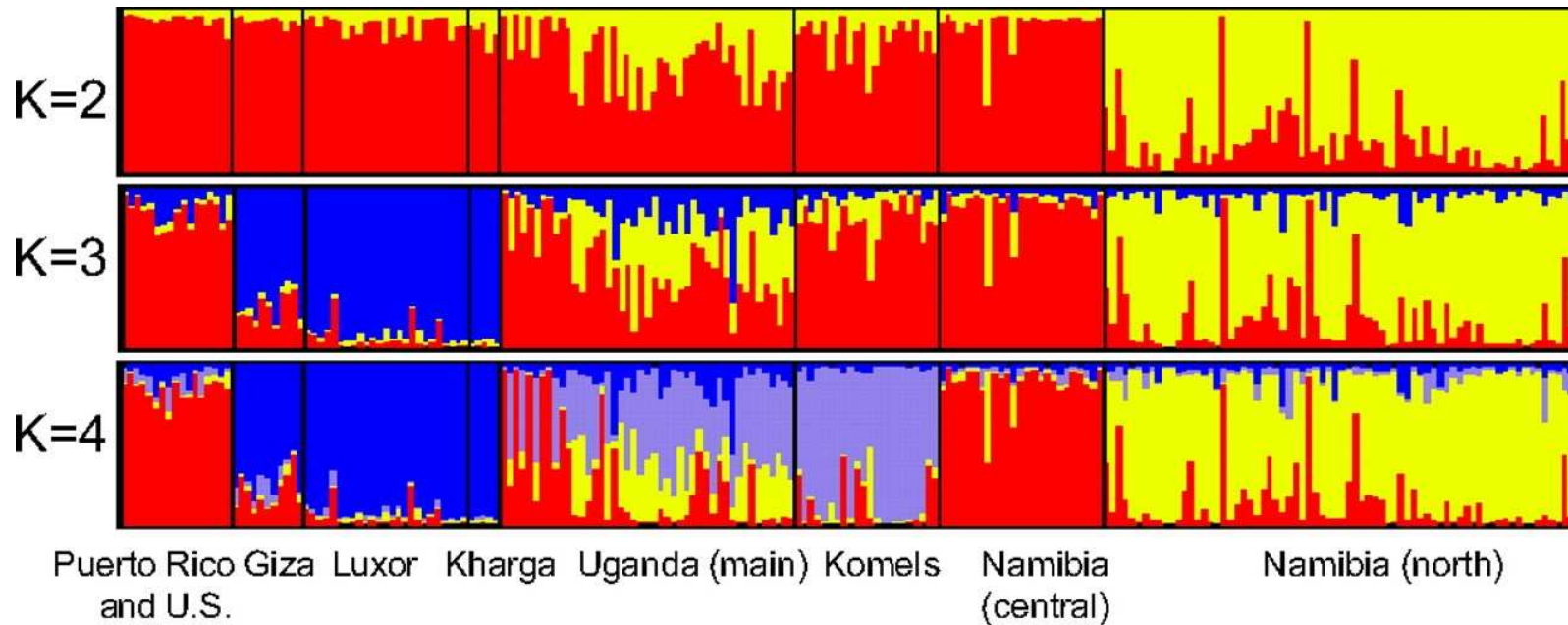
- Doteraz sme predpokladali, že nová generácia vzniká **náhodným párovaním** (random mating)
- Väčšina organizmov sa vyvíja v **subpopuláciách**, s obmedzeným prenosom genetického materiálu medzi subpopuláciami
- Frekvencie toho istého SNPu v dvoch subpopuláciách môžu byť značne odlišné
- \Rightarrow “falošné” korelácie medzi SNPami (napr. aj medzi chromozómami), ak pracujeme s viacerými subpopuláciami naraz
- \Rightarrow chybné výsledky pri LD a WGAS

Príklad: frekvencie alel jedného konkrétneho SNPu u ľudí v rôznych častiach sveta



zdroj: genome.ucsc.edu

Štruktúra populácie psov



Boyko et al. PNAS 2009; software STRUCTURE Pritchard et al. Genetics 2000

- Program STRUCTURE rozdelí populáciu na K subpopulácií (farby)
- Každý stĺpec je jedinec z populácie
- Pomer farieb zodpovedá pomeru SNPov z každej z K populácií

Ako funguje STRUCTURE?

- **Vstup:** Vzorka haplotypov X , ktorú chceme rozdeliť do K subpopulácií
- Definujeme stochastický model s nasledujúcimi premennými:
 - $P_{i,j}$ - frekvencia SNPu j v subpopulácii i
 - Q_i - aká časť SNPov v haplotype i patrí ku ktorej subpopulácii
 - $Z_{i,j}$ - priradenie subpopulácie SNPu j v haplotype i
- Model definuje $\Pr[X | P, Q, Z]$ a apriórne rozdelenie pre P, Q
- **Výstup:** $E[Q | X]$

Algoritmus Markov Chain Monte Carlo (MCMC)

- Premenné:
 - $P_{i,j}$ - frekvencia SNPu j v populácii i
 - $Z_{i,j}$ - priradenie subpopulácie SNPu j v haplotype i
 - Q_i - aká časť SNPov v haplotype i patrí ku ktorej populácii
- Začni s hodnotami $P^{(0)}, Z^{(0)}, Q^{(0)}$. V každej ďalšej iterácii získame novú náhodnú vzorku:
 - Vyber náhodnú vzorku $P^{(i)}, Q^{(i)}$ z distribúcie $\Pr(P, Q | X, Z^{(i-1)})$
 - Vyber náhodnú vzorku $Z^{(i)}$ z distribúcie $\Pr(Z | X, P^{(i)}, Q^{(i)})$
- Pre vhodné m, c , priemer postupnosti

$$Q^{(m)}, Q^{(m+c)}, Q^{(m+2c)}, \dots$$

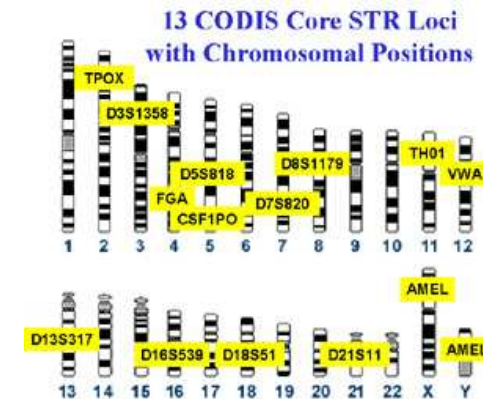
konverguje k hodnote $E[Q | X]$

Zhrnutie

- **SNPy (single nucleotide polymorphisms)** priebežne vznikajú a zanikajú v populáciách
- Ich frekvencia ovplyvnená navyše prirodzeným výberom
- Bez rekombinácie korelácia medzi SNPmi na tom istom chromozóme (**linkage disequilibrium**)
- Rekombinácie vytvárajú v genóme LD bloky
- Prítomnosť LD blokov možno využiť pri mapovaní asociácií znakov (**whole-genome association mapping**)
- Pravdepodobnostné modely veľkosti LD blokov, frekvencií alel, heterozygocity a pod. nám môžu veľa prezradiť o **histórii populácie**
- Pri analýzach treba brať do úvahy **štruktúru populácie**, ktorú možno odhadnúť pomocou výpočtových metód

Ďalšie typy polymorfizmov

- **Krátke indely**
- **Mikrosatelity a minisatelity** (jednoduché krátke opakujúce sa sekvencie)
13 lokusov ako štandardný “odtlačok” pre porovnávanie DNA vzoriek na súdoch v USA



- **Transpozóny** (Alu, LINE, SINE)
Alu má cca milión kópií, cca 1 nová kópia na 20 novorodencov
- **Veľké úseky s variabilnou multiplicitou** (Large scale copy number variations)