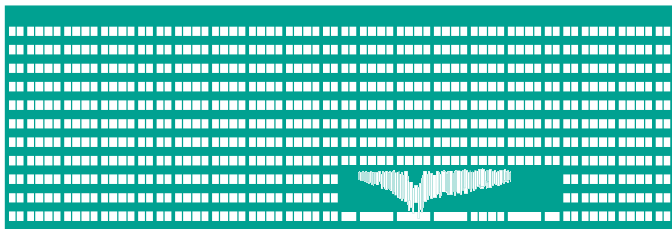


VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Algoritmy pro Bioinformatiku

Predikce genů

Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

6. listopadu 2019



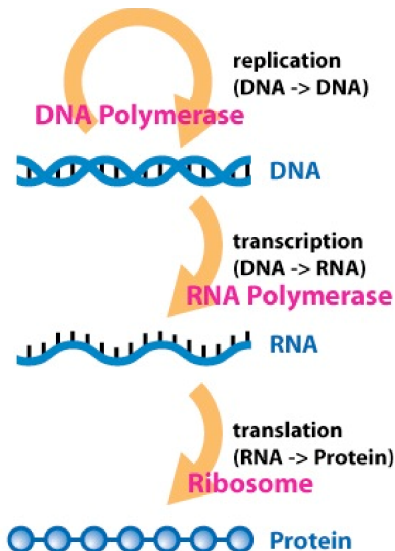
- Připomenutí pojmů
- Statistické metody
- Metody založené na homologii sekvencí
- Srovnávací metody



- V roce 1961 Sydney Brenner a Francis Crick objevili, že pokud dojde v DNA ke smazání jednoho či dvou nukleotidů, pak dojde k dramatické změně ve struktuře proteinu, tzv. frameshifting mutations. Pokud jsou však smazány tři po sobě jdoucí nukleotidy, pak je změna ve struktuře výsledného proteinu velmi malá.
- Důsledek: gen a jeho protein jsou kolineární strukturou s přímou korelací mezi trojicemi nukleotidů a aminokyselinami v proteinu.
- V roce 1977 překvapivé zjištění - gen není tvořen celou sekvencí, ale skládá se z několika podsekvencí oddělenými nekódujícími oblastmi => exony a introny.



- Původně vysloveno v roce 1958.
- V 60 letech napadeno Teminem a kol., nalezena možnost přenosu virové RNA do DNA.
- Transkripce a translace pouze genů.



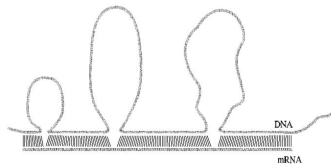
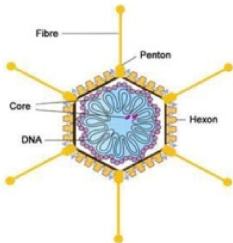


- Kodon - trojice nukleotidů.
- 64 kodonů je namapováno na 20 aminokyselin a STOP kodony.
- Některé aminokyseliny kódovány více kodony.
- Bodová mutace - substituce:
 - synonymní => nedojde ke změně proteinu.
 - nesynonymní => změna v aminokyselině vede k produkci jiného proteinu.

AMINOACID	CODONS	REDUNDANCY
ALANINE	GC*	4
CYSTEINE	TGC,TGT	2
ASPARTIC ACID	GAC,GAT	2
GLUTAMINE ACID	GAA,GAG	2
PHENYLALANINE	TTC,TTT	2
GLYCIN	GG*	4
HISTIDINE	CAC,CAT	2
ISOLEUCINE	ATA,ATC,ATT	3
LYSINE	AAA,AAG	2
LEUCINE	CT*,TTA,TTG	6
METHIONINE	ATG	1
ASPARGINE	AAC,AAT	2
PROLINE	CC*	4
GLUTAMINE	CAA,CAG	2
ARGININE	AGA,AGG,CG*	6
SERINE	AGC,AGT,TC*	6
THREONINE	AC*	4
VALINE	GT*	4
TRYPTOPHAN	TGG	1
TYROSINE	TAC,TAT	2
STOP	TAA,TAG,TGA	3



- Na základě experimentů s bakteriemi bylo zjištěno, že sekvence DNA, RNA a proteinů jsou kolineární. Předpokládalo se, že u eukaryotických buněk (organismů) to bude stejné.
- V roce 1977 Phillip Sharp a Richard Roberts experimentovali s mRNA virového proteinu, tzv. hexonu.
- Namapovali vzniklou mRNA zpět na DNA viru a výsledek zkoumali elektronovou mikroskopií.
- mRNA - DNA hybrid utvořil zvláštní tři smyčkovou strukturu namísto spojitého namapování.

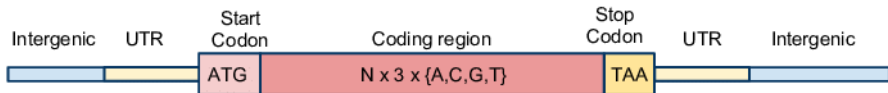




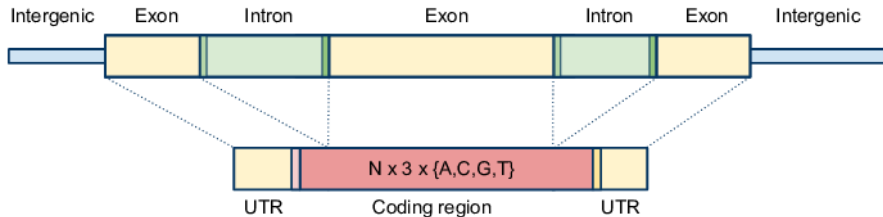
- Prokaryotické organismy nemají introny - jejich geny tvoří spojitě úseky DNA.
- U eukaryotických organismů gen je kombinací kódujících úseků (exonů) a nekódujících úseků (intronů).
- Eukaryotické geny jsou obtížněji predikovatelné.

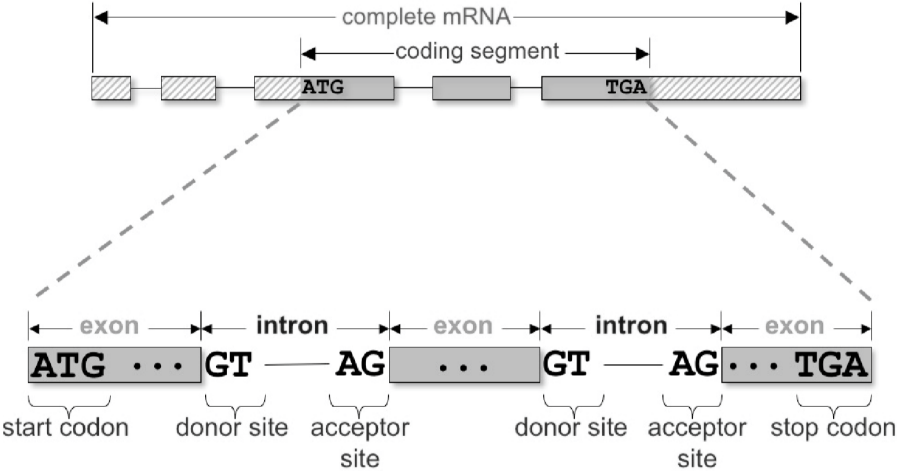


A) Prokaryotic Gene



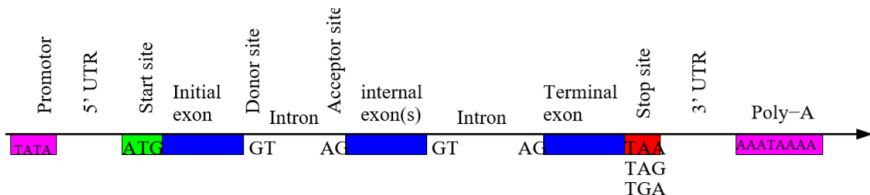
B) Eukaryotic Gene







- Exon - kódující oblast eukaryotického genu.
- Intron - nekódující oblast eukaryotického genu.
- Mezigenová oblast - oblasti DNA mimo úseky genů.
- UTR - untranslated terminal region, nekódující oblasti, lokalizovány před místem, kde dochází k inializaci translace a za místem ukončujícím translaci.
- Promoter - Oblast v mezigenové oblasti genomu, předchází gen a nepřekrývá se s ním. Iniciace transkripce.





- Osekvenována celá řada genomů.
- Zajímavé sekvence jsou sekvence kódující geny => zdrojové kódy proteinů.
- Nejspolehlivější metoda osekvenování mRNA => pracné a nákladné.
- Potřeba vývoje automatizovaných postupů pro predikci genů ze znalosti sekvencí genomu.



Predikce genů

Mějme sekvenci DNA. Problém predikce genů znamená nalezení všech genů, které jsou v dané sekvenci obsaženy.

- Sekvence kódující geny nejsou náhodnými sekvencemi.
- Pořadí kodonů splňuje určitá biologická pravidla a zachovává se v průběhu evoluce.
- Předpokládáme, že geny a mezigenové oblasti mají různé rozdělení nukleotidů.
- U eukaryot předpokládáme různé rozdělení kodonů v exonech a intronech.



- Přímé:
 - Sekvenování mRNA, nebo proteinových sekvencí.
- Výpočetní:
 - Hledání shody s již známými geny. Metody založené na tzv. homologii.
 - Hledání shody statistických vzorů společných všem genům. Tzv. metody ab initio.
 - Hybridní metody.



Deoxyribonukleová kyselina, běžně označovaná **DNA** (z anglického deoxyribonucleic acid, česky zřídka i **DNK**), je **nukleová kyselina**, nositelka **genetické informace** všech organismů s výjimkou některých **nebuněčných**, u nichž hraje tuto úlohu **RNA** (např. **RNA viry**). DNA je tedy pro **život** velmi důležitou látkou, která ve své struktuře kóduje a **buňkám** zadává jejich program a tím předurčuje **vývoj** a vlastnosti celého **organismu**. U **eukaryotických** organismů (jako např. **rostliny** a **živočichové**) je DNA hlavní složkou **chromatinu**, směsi nukleových kyselin a proteinů, a je uložena zejména uvnitř **buněčného jádra**, zatímco u **prokaryot** (např. **bakterie** a **archea**) se DNA nachází volně v **cytoplazmě**.

Oba texty jsou vytvořeny v českém jazyce nad stejnou abecedou, různé posloupnosti a četnosti slov nás však vedou k textům popisujícím různá témata.

Počítač je v **informatice** zařízení a **výpočetní technika**, která zpracovává **data** pomocí předem vytvořeného **programu**. Současný počítač je **elektronický** a skládá se z **hardwaru**, který představuje fyzické části počítače (**mikroprocesor**, **klávesnice**, **monitor** atd.) a ze **softwaru** (**operační systém** a **programy**). Počítač je zpravidla ovládán **uživatelem**, který poskytuje počítači data ke zpracování prostřednictvím jeho **vstupních zařízení** a počítač výsledky prezentuje pomocí **výstupních zařízení**. V současnosti jsou počítače využívány téměř ve všech oborech lidské činnosti.



- ORF
- Četnosti kodonů.
- Četnosti hexamerů - oligonukleotidy o délce 6 nt.
- Promotery genů.
- Místa přechodu exon/intron (splice sites).
- UTR - untranslated terminal regions.



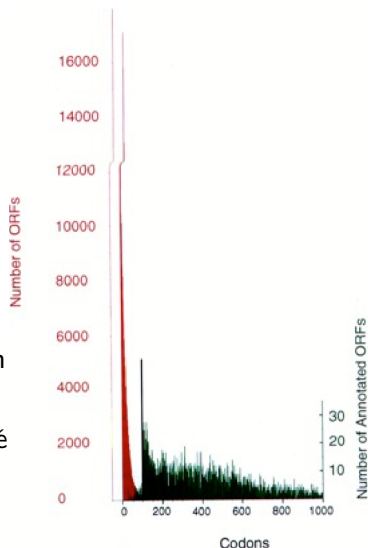
Open Reading Frame

Open Reading Frame (ORF) je jakákoli sekvence začínající start kodonem, končící end kodonem a neobsahující žádný end kodon uvnitř sekvence.

- Účel: detekovat potenciální kódující úseky.
- Každá DNA sekvence má šest možných ORF => předpokládáme, že kodony (trojice nukleotidů) čteme od pozice i , $i + 1$ a $i + 2$ v obou směrech: na přímém i komplementárním vláknu.
- Ne každý ORF kóduje gen.
- Velikost je dělitelná třemi.
- Odhalení rozdílu mezi genem a náhodným ORF je klíčovou úlohou.



- Základem je hledat ORF, které vypadají jako geny.
- Počet všech ORF je mnohokrát vyšší než počet anotovaných ORF.
- Hledáme delší ORF, cca 100 nt a více.
- Předchozí bod platný především u prokaryot.
- Krátké ORF jsou poměrně časté u eukaryotických organismů.





- Pokud by DNA byla náhodnou sekvencí stejně pravděpodobných nukleotidů, pak bychom očekávali, že délka jednoho genu bude cca 21 nukleotidů.
- Obecně u bakterií jsou geny poměrně dlouhé, stovky nukleotidů.
- Počet kodonů v průměrném proteinu bakterie je přibližně 300.
- Průměrná délka genu u obratlovců 30kb, průměrná délka exonu 1-2kb.
- Prokaryoty nemají introny a mezigenové oblasti jsou krátké.
- Geny se navíc mohou překrývat nebo jsou dokonce vnořené v sobě.
- Místa začátku translace obtížně určitelná.



- ORF se mohou překrývat.
- Pokud stanovíme určitou hranici délky genu nebudeme schopni detekovat krátké geny.
- U prokaryot se vyskytují alternativní start kodony, kromě ATG také GTG a TTG.
- ORF samotné nepostačují pro detekci genu, využívají se další charakteristiky.



- Rozdělení nukleotidů v okolí nějakého signálního místa (start kodonu, end kodonu, počátek a konec intronu).
- Používají se k odlišení konkrétních signálních oblastí od nesignálních oblastí.
- Například nukleotidy v okolí počátečního start kodonu:

Pos.	-8	-7	-6	-5	-4	-3	-2	-1	+1	+2	+3	+4	+5	+6	+7
A	.16	.29	.20	.25	.22	.66	.27	.15	1	0	0	.28	.24	.11	.26
C	.48	.31	.21	.33	.56	.05	.50	.58	0	0	0	.16	.29	.24	.40
G	.18	.16	.46	.21	.17	.27	.12	.22	0	0	1	.48	.20	.45	.21
T	.19	.24	.14	.21	.06	.02	.11	.05	0	1	0	.09	.26	.21	.21



- Měří relativní četnost určitého kodonu v daném organismu, genu nebo ORF.
- Například sekvence HIV-1 polymerazy: čtyři kodony kódující Alanine:

Kodon	Počet
GCA	41576
GCC	9461
GCG	1017
GCT	11031

- Vidíme značné zešíkmění rozdělení kodonů.



- Relativní použití synonymních kodonů (RSCU) pro dvojic (i, j) , kde i je aminokyselina a j je jeden z n_i synonymních kodonů. X_{ij} je počet j -tého kodonu aminokyseliny i .
- $RSCU > 1$ indikuje preferovaný kodon, zatímco $RSCU < 1$ nepreferovaný kodon.

$$RSCU_{ij} = \frac{X_{ij}}{\frac{1}{n_i} \sum_{k=1}^{n_i} X_k}$$

- Relativní adaptivnost:

$$w_{ij} = \frac{RSCU_{ij}}{RSCU_{max}} = \frac{X_{ij}}{max_j X_{ij}}$$



Kodon	Počet	RSCU	w
GCA	41576	2.64	1
GCC	9461	0.60	0.23
GCG	1017	0.064	0.02
GCT	11031	0.70	0.27



		<u>E.coli</u>		Yeast		<u>E.coli</u>		Yeast			
		RSCU	w	RSCU	w	RSCU	w	RSCU	w		
Phe	UUU	0.456	0.296	0.203	0.113	Ser	UCU	2.571	1.000	3.359	1.000
	UUC	1.544	1.000	1.797	1.000		UCC	1.912	0.744	2.327	0.693
Leu	UUA	0.106	0.020	0.601	0.117	UCA	0.198	0.077	0.122	0.036	
	UUG	0.106	0.020	5.141	1.000	UCG	0.044	0.017	0.017	0.005	
Leu	CUU	0.225	0.042	0.029	0.006	Pro	CCU	0.231	0.070	0.179	0.047
	CUC	0.198	0.037	0.014	0.003		CCC	0.038	0.012	0.036	0.009
	CUA	0.040	0.007	0.200	0.039		CCA	0.442	0.135	3.776	1.000
	CUG	5.326	1.000	0.014	0.003		CCG	3.288	1.000	0.009	0.002
Ile	AUU	0.466	0.185	1.352	0.823	Thr	ACU	1.804	0.965	1.899	0.921
	AUC	2.525	1.000	1.643	1.000		ACC	1.870	1.000	2.063	1.000
	AUA	0.008	0.003	0.005	0.003		ACA	0.141	0.076	0.025	0.012
Met	AUG	1.000	1.000	1.000	1.000	ACG	0.185	0.099	0.013	0.006	
Val	GUU	2.244	1.000	2.161	1.000	Ala	GCU	1.877	1.000	3.005	1.000
	GUC	0.148	0.066	1.796	0.831		GCC	0.228	0.122	0.948	0.316
	GUA	1.111	0.495	0.004	0.002		GCA	1.099	0.586	0.044	0.015
	GUG	0.496	0.221	0.039	0.018		GCG	0.796	0.424	0.004	0.001



- Pozorovaný CAI_{obs} sekvence o L kodonech je dán geometrickým průměrem každého z kodonů:

$$CAI_{obs} = \left(\prod_{k=1}^L RSCU_k \right)^{1/L}$$

- Tato hodnota je poté porovnána s maximálním možným CAI, který vytvoříme ze všech možných sekvencí kodonů stejné délky, které kódují stejný protein:

$$CAI_{max} = \left(\prod_{k=1}^L RSCU_{max} \right)^{1/L}$$

$$CAI = \frac{CAI_{obs}}{CAI_{max}}$$

Codon Usage in Human Genome

	U	C	A	G
U	UUU Phe 57	UCU Ser 16	UAU Tyr 58	UGU Cys 45
	UUC Phe 43	UCC Ser 15	UAC Tyr 42	UGC Cys 55
	UUA Leu 13	UCA Ser 13	UAA Stp 62	UGA Stp 30
	UUG Leu 13	UCG Ser 15	UAG Stp 8	UGG Trp 100
C	CUU Leu 11	CCU Pro 17	CAU His 57	CGU Arg 37
	CUC Leu 10	CCC Pro 17	CAC His 43	CGC Arg 38
	CUA Leu 4	CCA Pro 20	CAA Gln 45	CGA Arg 7
	CUG Leu 49	CCG Pro 51	CAG Gln 66	CGG Arg 10
A	AUU Ile 50	ACU Thr 18	AAU Asn 46	AGU Ser 15
	AUC Ile 41	ACC Thr 42	AAC Asn 54	AGC Ser 26
	AUA Ile 9	ACA Thr 15	AAA Lys 75	AGA Arg 5
	AUG Met 100	ACG Thr 26	AAG Lys 25	AGG Arg 3
G	GUU Val 27	GCU Ala 17	GAU Asp 63	GGU Gly 34
	GUC Val 21	GCC Ala 27	GAC Asp 37	GGC Gly 39
	GUA Val 16	GCA Ala 22	GAA Glu 68	GGA Gly 12
	GUG Val 36	GCG Ala 34	GAG Glu 32	GGG Gly 15

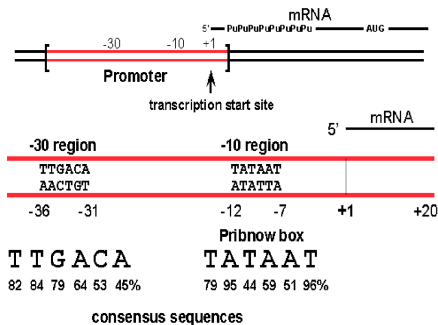
10.



- Začátky a konce exonu jsou signalizovány dvojicemi dvou-nukleotidových úseků GT a AC tzv. donor a acceptor sites.
- GT a AC jsou krátké sekvence a vyskytují se velmi často => detekce je obtížná.
- Vykazují však zajímavé statistické okolí.



Promoter structure in prokaryotes



- Transkripce začíná indexem 0.
- Pribnow Box (-10)
- Gilbert Box (-30)
- Místo nasednutí ribozomu (+10)



- Charakteristiky použity k vytvoření trénovacích datasetů.
- Detekce se stává klasifikačním problémem.

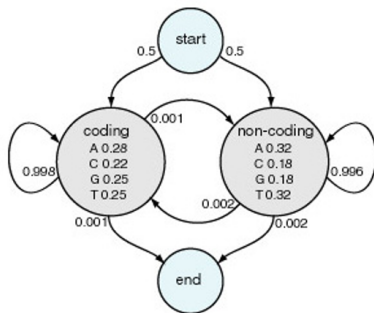
Table 2. Percentage accuracy (average of specificity and sensitivity) of the coding measures in predicting region coding

Measure	Human 54 Penrose	Human 108 Penrose	Human 162 Penrose	<i>E. coli</i> 54 Penrose	Human 54 Classical
Hexamer	70.5	73.1	74.2	67.5	–
Position Asymmetry	70.2	76.6	80.6	61.6	70.3
Dicodon Usage	70.2	72.9	73.9	67.5	–
Fourier	69.9	76.5	80.8	61.3	69.9
Hexamer-1	69.9	72.6	73.8	66.8	–
Hexamer-2	69.9	72.6	73.8	66.7	–
Run	66.6	70.3	71.3	63.6	67.9
Codon Usage	65.2	68.0	69.5	64.1	66.
Repeat	65.1	69.9	73.1	62.4	–
Autocorrelation	64.9	71.1	77.0	58.2	64.9
Dinucleotide Bias	62.9	55.5	50.7	55.9	62.7
Diamino Acid Usage	62.8	66.3	67.8	61.3	–
Composition	61.7	64.1	65.9	61.7	61.3
Amino Acid Usage	60.6	63.4	64.7	59.7	61.3
Word	59.5	66.4	71.4	56.6	61.0
Entropy	58.4	63.1	66.2	55.0	58.4
Dinucleotide Frame	58.0	62.9	66.6	54.6	58.0
Open Reading Frame	57.8	59.2	60.7	57.4	57.8
Stability Hydrophobicity	55.5	57.5	58.7	55.5	55.5
Codon Prototype	54.7	56.1	56.4	54.7	54.7
Period 9	52.5	53.0	52.8	51.8	52.4

Data from five benchmark situations are shown, with varying data set (Human or *E. coli*), window length (54, 108, or 162) and decision method (Penrose discriminant or Classical linear discriminant).



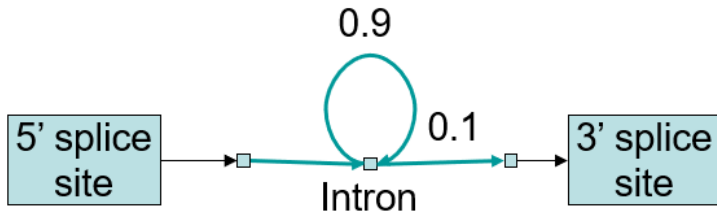
- První program založený na Markovových modelech GeneMark.
- Předpoklad, že dlouhé ORF sekvence jsou geny => slouží jako trénovací dataset, např. Glimmer
- Obecně použijeme Markovův model pro ohodnocení pravděpodobnosti, že daná sekvence odpovídá modelu.
- Vybíráme takovou interpretaci sekvence, která má nejvyšší pravděpodobnost.
- Markovův model sestavujeme z již existujících anotovaných sekvencí genů.



$$\Phi = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0.5 & 0.998 & 0.002 & 0 \\ 0.5 & 0.001 & 0.996 & 0 \\ 0 & 0.001 & 0.002 & 0 \end{bmatrix}$$

$$H = \begin{bmatrix} 0.28 & 0.32 \\ 0.22 & 0.18 \\ 0.25 & 0.18 \\ 0.25 & 0.32 \end{bmatrix}$$

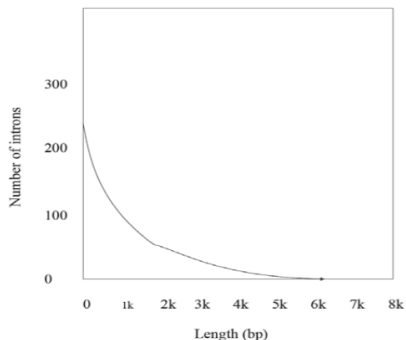
- $x_m(i)$ - pravděpodobnost, že jsme ve stavu m na pozici i
- $H(m, y_i)$ - pravděpodobnost, že zvolíme symbol y_i ve stavu m .
- Φ_{mk} - Pravděpodobnost, že přejdeme ze stavu k do stavu m .



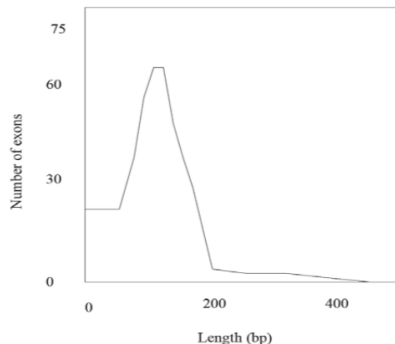
- Předpokládejme, že každá modrá šipka vyprodukuje jeden symbol.
- Jaká je pravděpodobnost, že intron bude mít délku právě 4 nukleotidy (uvažujeme pouze GT,AG)? 10%
- Jaká je pravděpodobnost, že intron bude mít délku právě 5 nukleotidů (uvažujeme GT,AG a jeden nukleotid)? 9%
- Jaká je pravděpodobnost, že intron bude mít délku právě 6 nukleotidů (uvažujeme GT,AG a dva nukleotidy)? 8.1%



Výpočet délek Markovových řetězců pro exony/introny vede na geometrické rozdělení: $p^k(1-p) \Rightarrow$ neodpovídá realitě:



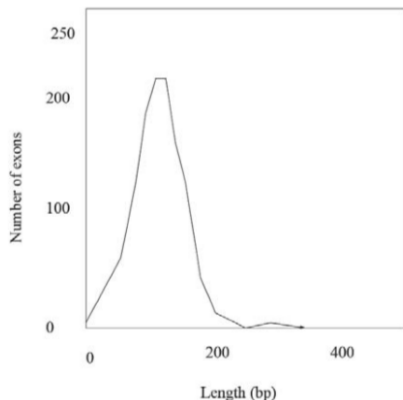
introns



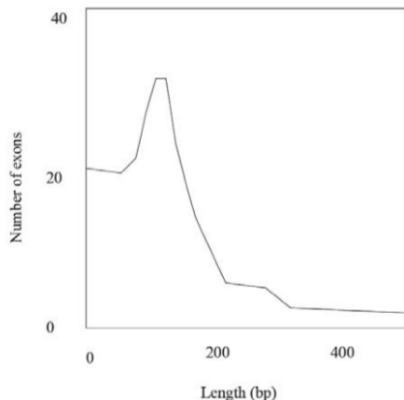
initial exons



Výpočet délek Markovových řetězců pro exony/introny vede na geometrické rozdělení: $p^k(1-p) \Rightarrow$ neodpovídá realitě:



internal exons



terminal exons



- Například nástroj Prodigal.
- Nepoužívá dlouhé ORF sekvence => vedou k chybným predikcím v případech sekvencí bohatých na GC nukleotidy.
- Stop kodony: TAA, TGA a TAG obsahují AT, jejich četnost u organismů s vysokým obsahem GC nukleotidů je nižší a tudíž vede k vyšší pravděpodobnosti zachycení delších ORF.
- V kódující DNA obsah GC na třetí pozici v kodonu je vyšší oproti nekódující DNA.



- Evoluce kódujících oblastí je pomalejší, než evoluce nekódujících oblastí.
- Umožňuje porovnat genomovou sekvenci s příbuzným organismem a najít podobné anotované sekvence.
- Nejrozšířenější nástroj BLAST.
- Nejspolehlivější metoda pro predikci genů.
- Společně se sekvencí obvykle obdržíme také anotaci (k určení funkce genu).
- Limitem metody je neúplný obsah databází proteinových sekvencí. Některé příbuzné organismy stále nemusí být vůbec v databázích.



Daná DNA sekvence je porovnána s podobnou DNA sekvencí z evolučně blízkého organismu. Geny jsou predikovány v obou sekvencích na základě hypotézy, že exony se při evoluci zachovávají, zatímco introny ne.

- CEM - conserved exon method.
- TWINSKAN - použití Markovových modelů a podobnosti genomů..



- Obecně jsou softwarové nástroje specifické pro každý organismus.
- Nejlépe fungují na genech, které jsou podobné nějakým již anotovaným genům.
- Schopny nalézt protein kódující oblasti daleko lépe než nekódující oblasti.
- Všechny modely jsou nedokonalé.



- Stále spousta prostoru pro další rozvoj => především porozumění inicializace translace.
- Schopnost predikovat geny je silně závislá na GC obsahu sekvencí.
- Odhadovalo se, že 10-30% anotovaným genů v databázích jsou náhodné ORF.
- Shrnutí: neexistuje nástroj, který by se 100% přesností našel všechny geny.

DĚKUJI za pozornost

Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

6. listopadu 2019