

Predikce genů

Pro zajímavost, nebude součástí zkoušky...

Důležité, pravděpodobně bude u zkoušky...

Molekulárně biologická data

- **Výkonné technologie:**

Automatické sekvenování

MALDI-TOF

NMR spektroskopie

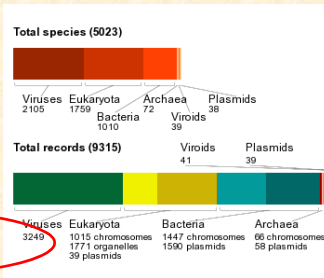
Proteinová krystalografie

Výrazný nárůst množství biologických dat.

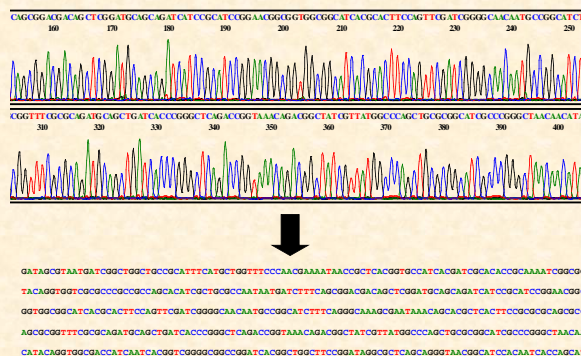
Rozdělení molekulárně biologických databází

- **Databáze:**
- Primární
- Sekundární
- Strukturální

Genomové zdroje



Molekulárně biologická data



„Syrové“ sekvence DNA



Identifikace a anotace genů a proteinů

Table 1
Software commonly used for bacterial genome annotation and comparison

| Software | Description |
|----------------------------------|-------------------------------------------------------------------------------------------------------------------------------|
| DNA level annotation | |
| GeneMark | http://www.gatech.edu/genemark/ |
| Glennier | http://www.genomics.jhu.edu/Glennier/ |
| SHOW | http://genome.jgi-psf.org/SHOW/ |
| rNAscan-SE | http://www.fwlab.ucsc.edu/rNAscan-SE/ |
| rNAmmer | http://www.cbs.dtu.dk/services/rNAmmer/ |
| RapSearch | http://www.ncbi.nlm.nih.gov/COG/databases/RapSearch/ |
| IslandPath | http://www.pathogenomics.sfu.ca/islandpath/ |
| Protein level annotation | |
| BLAST | http://www.ncbi.nlm.nih.gov/BLAST/ |
| InterProScan | http://www.ebi.ac.uk/interproscan/ |
| COGNITOR | http://www.ncbi.nlm.nih.gov/COG/databases/cognitor.html |
| PRISM | http://bioinfo.genopole-toulouse.prd.fr/prim/ |
| GOAnno | http://pips.u-strasbg.fr/GOAnno/ |
| PSORTb | http://www.psort.org/psortb/ |
| TMDMM | http://www.cbs.dtu.dk/services/TMDMM/ |
| SignalP | http://www.cbs.dtu.dk/services/SignalP/ |
| Comparative genomic tools | |
| MaVe | http://lar-lab.usc.edu/mauve/ |
| MOSMIC | http://img.jgi-psf.org/mauve/mauve.html |
| ACT | http://www.sanger.ac.uk/software/ACT/ |
| CLUST | http://img.jgi-psf.org/mauve/mauve.html |
| MaGe | http://www.genoscope.cn.fr/ago/mauve/ |
| Pathologic | http://bioeye.org/ |
| PFMA2 | http://compbio.ucsf.edu/pfma2/ |
| The SIFTS | http://www.ebi.ac.uk/sift/ |
| STRING | http://string.embl.de/ |
| PyPh | http://www.cbs.dtu.dk/services/pyph/ |
| HoSeq | http://pbl.univ-lyon1.fr/software/HoSeq/ |

Predikce genů kódujících proteiny

- **Prokaryotické geny**
- Nepřerušované úseky DNA mezi **startovním kodonem** (ATG, GTG, TTG, CTG) a **stop kodonem** (TAA, TGA, TAG).
- **Eukaryotické geny**
- Přerušovány **introny**. Průměrná délka exonu je 50 kodonů, některé jsou mnohem kratší.
- Některé introny extrémně dlouhé, geny zabírají mbp v genomové DNA.

Predikce eukaryotických genů je mnohem složitější než predikce genů prokaryotických a představuje **STÁLE NEVYŘEŠENÝ** problém!

Prokaryotické geny

- **Prokaryotický gen = nejdelší ORF odpovídající danému úseku DNA.**

```
GTATGCTGGTATTGTGGATGCCGTTACCCCTGCTGAGCCCTATCGGAAAGCCAGCGTGTATCCGGCCGCC
GACCGTATTGATGGTCGCCAAGCTGATGTTGTTAGCCGGGCGATGCCCGCAGCTGGCCATAACGATAGC
CGTCTGTTTACCGGTCTGAGCCCGGGTATCAGCTGCATCTGCGGAAACCGCGCTGCGCTGCGGGCGAAG
TGAGCGTCTGTTTATCGCTTTCGCTTGCCTGAAAGATCGCGCATTTGTCGCCGATCGAACTGGAATGGGTG
TGCGCCACCGCGCTTCGGATGCGGATGATCTGCTGCATCCGAGCTGTCGTCGCGTGAAGATCATTATTGG
CGCAGCATGTGCTGGCGCGCGCCAGCCACCTGTACCGCCGATTTGGCGTGTGCGATCGTGATGGCACCG
TGACCGTATTATTTGTTGGGAAACCGAGCATGAAATGCGGGCAGCCAGCCGGAACCAACAGCCGAGCTT
TAAACCGAGCAGCATCCGATGCGCACTTTAGCGTCCCGCGATACCCGCTTAAAGGATCTTCTATCGG
ACCGCGCGATCTGAGGATCGAAACTGTTATGATGATGCGCGGAAACCGCGCCGACCTTTGTGGGTA
ACAGCGAAGATGGTGTGCGTCTGTTTACCCCTGAATAGCAAAGTGTGTAATAATTCGATTGAAAGCGAGCCGAA
CGCCCTCAGAGCGGACGATGCCCTCTGCGCGCCGCTGAGCCCGGCGATACCGTGTGCGTGGCTGGCTGGC
GCCCGAAGATGGTCCGATGCGGATTAATGATGGCATTGTTATCTGCAATGGCCGATTAACCTAATGGG
```

Překlad DNA sekvence

- **ExpASY**
<http://www.expasy.org/tools/dna.html>
- **ORF Finder (NCBI)**
<http://www.ncbi.nlm.nih.gov/gorf/gorf.html>

nonpolar polar basic acidic (stop codon)

Překlad DNA sekvence

The table shows the 64 codons and the amino acids for each. The direction of the mRNA is 5' to 3'.

| 1st base | 2nd base | | | |
|----------|----------------------------------------------|---------------------------|---------------------------|------------------------|
| | U | C | A | G |
| U | UUU (Phe/F) Phenylalanine | UCU (Ser/S) Serine | UAU (Tyr/Y) Tyrosine | UGU (Cys/C) Cysteine |
| | UUC (Phe/F) Phenylalanine | UCC (Ser/S) Serine | UAC (Tyr/Y) Tyrosine | UGC (Cys/C) Cysteine |
| | UUA (Leu/L) Leucine | UCA (Ser/S) Serine | UAA (Ochne/S) Stop | UGA (Ogal/S) Stop |
| | UUG (Leu/L) Leucine | UCG (Ser/S) Serine | UAG (Amber/S) Stop | UGG (Trp/W) Tryptophan |
| C | CUU (Leu/L) Leucine | CCU (Pro/P) Proline | CAU (His/H) Histidine | CGU (Arg/R) Arginine |
| | CUC (Leu/L) Leucine | CCC (Pro/P) Proline | CAC (His/H) Histidine | CCG (Arg/R) Arginine |
| | CUA (Leu/L) Leucine | CCA (Pro/P) Proline | CAA (Gln/Q) Glutamine | CCA (Arg/R) Arginine |
| | CUG (Leu/L) Leucine | CCG (Pro/P) Proline | CAG (Gln/Q) Glutamine | CGC (Arg/R) Arginine |
| A | AUU (Ile/I) Isoleucine | AUU (Thr/T) Threonine | AUU (Asn/N) Asparagine | AUU (Ser/S) Serine |
| | AUC (Ile/I) Isoleucine | AUC (Ile/I) Isoleucine | AAC (Asn/N) Asparagine | ACC (Ser/S) Serine |
| | AUA (Ile/I) Isoleucine | AUA (Thr/T) Threonine | AUA (Lys/K) Lysine | AUA (Arg/R) Arginine |
| | AUG (Met/M) Methionine, Start ^(S) | AUG (Thr/T) Threonine | AAG (Lys/K) Lysine | AAC (Arg/R) Arginine |
| G | GUU (Val/V) Valine | GUU (Asp/D) Aspartic acid | GAU (Asp/D) Aspartic acid | GGU (Gly/G) Glycine |
| | GUC (Val/V) Valine | GUC (Ala/A) Alanine | GAC (Asp/D) Aspartic acid | GGU (Gly/G) Glycine |
| | GUA (Val/V) Valine | GUA (Ala/A) Alanine | GAA (Glu/E) Glutamic acid | GGA (Gly/G) Glycine |
| | GUG (Val/V) Valine | GUG (Asp/D) Aspartic acid | GAG (Glu/E) Glutamic acid | GGG (Gly/G) Glycine |

ExpASY

<http://www.expasy.org/tools/dna.html>

Site Map Search ExpASY Contact us

Search Swiss-Prot/TrEMBL for Go Clear

The ExpASY Server requires Javascript to be fully functional. You may not see all the information available for this page ([More information](#)).

ExpASY Proteomics Server

The ExpASY (Expert Protein Analysis System) proteomics server of the Swiss Institute of Bioinformatics (SIB) is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE. (Disclaimer / References / Linking to ExpASY)

[Databases] [Tools & Software] [Education & Services] [Links] [Announcements] [Mirror Sites] [Job openings]

| Databases | Tools and software packages |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"> • UniProt Knowledgebase (Swiss-Prot and TrEMBL) - Protein knowledgebase • ViralZone - Portal to viral UniProtKB/Swiss-Prot entries *** • PROSITE - Protein families and domains • SWISS-2DPAGE - Two-dimensional polyacrylamide gel electrophoresis • World-2DPAGE Repository - A public standards-compliant repository for gel-based proteomics data published in the literature • MAPE-DB - A public repository for MAPE Gel electrophoresis documents | <ul style="list-style-type: none"> • Proteomics and sequence analysis tools <ul style="list-style-type: none"> ◦ Identification and characterization (Aldester, FindMod, Popitam, Phenix, MW, ProtParam...) ◦ DNA → Protein ◦ Similarity search (BLAST...) ◦ Protein profile searches (ScanProsite...) ◦ Post-translational modification and topology prediction ◦ Primary structure analysis ◦ Secondary and tertiary structure tools (Swiss-PdbViewer...) ◦ Alignment and Phylogenetic analysis • Melanin / ImageMaster - Software for 2-D PAGE analysis |

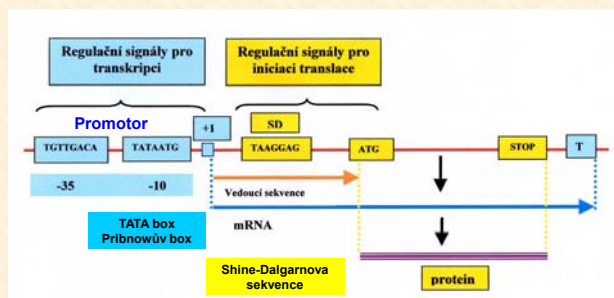
Prokaryotické geny

- **Velmi jednoduchý přístup k predikci genů**
Zjednodušení vede k chybám, ale jejich množství je **POMĚRNĚ MALÉ**.
- **Chyby mohou vznikat při SEKVENCOVÁNÍ DNA.**
Přidání/odstranění startovního a/nebo stop kodonu může vést ke **ZKRÁCENÍ**, **PRODLOUŽENÍ** nebo úplnému **VYNECHÁNÍ** genu.

Opravdu ORF kóduje protein?

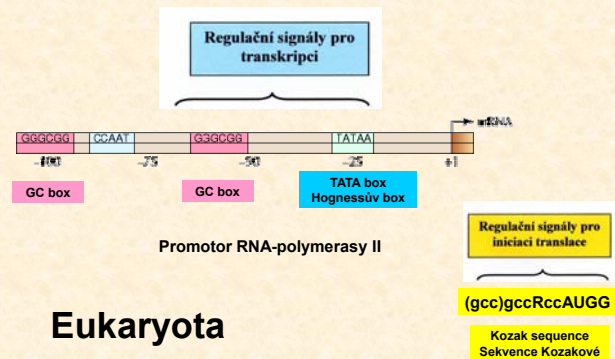
- **ORF kóduje protein, který je podobný již dříve popsanému proteinu** (prohledávání DATABÁŽÍ pomocí ALIGNMENTU).
- **ORF má typický obsah GC nebo frekvenci kodonů.** Srovnání s charakteristickými vlastnostmi známých genů ze stejného organismu.
- **Před ORF se nachází typické RBS (ribosome-binding site) nebo promotor.**

Translační a transkripční signální sekvence



Prokaryota

Translační a transkripční signální sekvence



Eukaryota

Opravdu ORF kóduje protein?

- **ORF kóduje protein, který je podobný již dříve popsanému proteinu** (prohledávání DATABÁŽÍ pomocí ALIGNMENTU) = **nejspolehlivější ověření**.
- **Nástroje pro překlad DNA jsou propojeny s prohledáváním databází.**

ORF Finder (NCBI)

<http://www.ncbi.nlm.nih.gov/gorf/gorf.html>

View 1 GenBank Redraw 100 SixFrames

| Frame | from | to | Length |
|-------|-------|-----|--------|
| +3 | 3.872 | 870 | |
| -2 | 1.857 | 857 | |

23 Frame 3
MetLVVDAVTLSSAYPEASRDPAAPTVIDGRHLYVSPGDAAGLGHNSRLFTGLSPGDDLHRETALALRAEVSVLFIRFALKDAGIVAPI
ELEVRDAAATAVPDAADLLHPSORPKDHWVRSVLAAGATTDTADFVCDRDGTSGYFRWETSIEIAGSSPDTKQDFKPSDRIGN
PSLPPNTAFKAIFYANAADRQDLKLFIDDAPEPAATFVGNSEDGVRLFTLNSGKIRIEASANGRSATDARLAPLSAGDTYVWLWGA
EDGADADYNDGIVLQWPI Stop

25 Frame 2
PLGNRPLONNNAIIIRIGTIFRAQPAQPHGIARADRRGTGIGRALTAVRARFNNTFTFAIQGKQHTIFAVTHKGGORFRRIINKQFQILT
RRVRIEDRFKQGIIRRDQAKVAIAARFNALFGIRLAARNFNAGFPFKITAHGAIITIAHRKIGOTGGRARRQHIAAPIMetIFQRTTARMetGGII
RIRNDDGITHFQDRGNAGIFGQKANKQHAHFRAGRGRGFAQMetLITRAQTGKQTAIVMetADLRGIARANNIVQATINHRRGGRTA
DFRIGAGGGGIHNGH

ORF Finder (NCBI)
<http://www.ncbi.nlm.nih.gov/gorf/gorf.html>

Eukaryotické geny Jednobuněčná eukaryota

- **Genomy jednobuněčných eukaryot se výrazně liší** (frekvence intronů, jak velká část genomu je tvořena geny kódujícími proteiny).
- *Saccharomyces cerevisiae* – 67% genomu je protein-kódující, jen 4% obsahují introny.
- Hlenky – průměrný gen obsahuje 3,7 intronu.
- Pro některá jednobuněčná eukaryota (kvasinky) je možné použít stejné postupy jako pro prokaryota.



Eukaryotické geny Mnohobuněčná eukaryota

- **Mnohobuněčná eukaryota**

Komplexní organizace genomu, geny separovány dlouhými **INTERGENOVÝMI** úseky, geny obsahují množství **INTRONŮ**, i velmi **DLOUHÝCH**.

Glyceraldehyd-3-fosfát-dehydrogenasa
Candida albicans

Eukaryotické geny Mnohobuněčná eukaryota

- **Mnohobuněčná eukaryota**

Komplexní organizace genomu, geny separovány dlouhými **INTERGENOVÝMI** úseky, geny obsahují množství **INTRONŮ**, i velmi **DLOUHÝCH**.

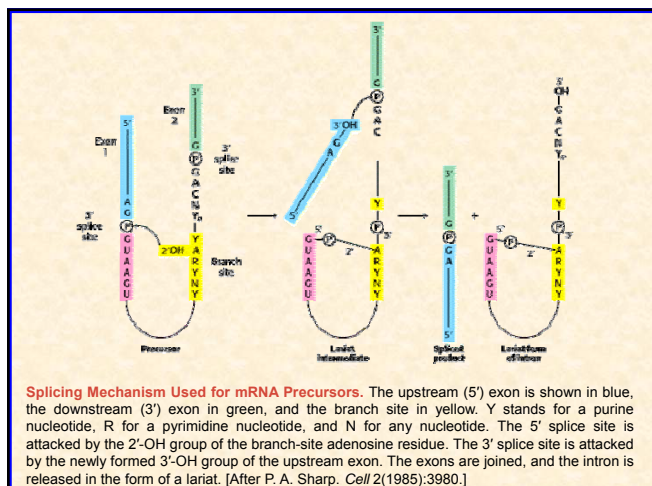
Glyceraldehyd-3-fosfát-dehydrogenasa
Homo sapiens

Eukaryotické geny Mnohobuněčná eukaryota

- **Rozpoznání exonů/intronů**

Identifikace míst sestřihu: **GT** na 5' konci, **AG** na 3' konci.
- **Chyby při rozpoznávání exonů/intronů**

Velké množství chyb. Dlouhé introny – určeny jako intergenové úseky. Krátké intergenové úseky – určeny jako introny.



Algoritmy a nástroje pro identifikaci genů

- **Predikce genů na základě sekvenční homologie** – vyhledávání v databázích pomocí algoritmů.
- **Predikce genů *ab initio*** – predikce na základě statistických parametrů DNA sekvence.
- **Většina běžně používaných metod kombinuje oba dva přístupy.**

Prokaryota

ATG.....TAA
Bez intronů
SEKVENČNÍ HOMOLOGIE

IDENTIFIKOVANÉ GENY VYUŽITY PRO „TRÉNOVÁNÍ“ STATISTICKÉ METODY

ANALÝZA ZBÝVAJÍCÍCH ČÁSTÍ GENOMU

Eukaryota

Mnoho intronů, dlouhé intergenové úseky
***Ab initio* STATISTICKÉ METODY**

IDENTIFIKOVANÉ EXONY

SEKVENČNÍ HOMOLOGIE

Algoritmy a nástroje pro identifikaci genů

- **Každý program má výhody a nevýhody – rozumné použít více predikčních nástrojů.**

GeneMark
GlimmerM
GRAIL
GenScan
Fgenes

Algoritmy a nástroje pro identifikaci genů

- **GeneMark**
<http://exon.gatech.edu/GeneMark>

Využívá **Markovovy** modely

Vyžaduje parametry specifické pro daný organismus = nutné „natrénování“ pomocí známých genů

Varianty pro prokaryotické, eukaryotické, virové sekvence

