

BIOINFORMATIKA V PRAXI

CVIČENÍ 3 – DRUHÁ ČÁST

IDENTIFIKACE GENŮ, PROTEINŮ A JEJICH FUNKCE

STUDIJNÍ MATERIÁLY

Studijní materiály předmětu C2130 Úvod do chemoinformatiky a bioinformatiky, přednáška **Predikce genu, Sequence-evolution-function: Computational Approaches in Comparative Genomics**.

CHYBY PŘI PREDIKCI GENŮ

Velmi zjednodušený přístup k predikci prokaryotických genů (genem je nejdelší ORF) vede k chybám, ale jejich množství je poměrně malé. Chyby mohou také vznikat při sekvencování DNA. Přidání/odstranění startovního a/nebo stop kodonu může vést ke zkrácení, prodloužení nebo úplnému vynechání genu.

ÚKOL 1 – příklady chyb vzniklých při sekvenaci

Pomocí predikčního programu **GeneMark** (<http://exon.gatech.edu/GeneMark>) identifikujte geny v sekvenačních výstupech a porovnejte je s původní sekvencí z databáze (část genomu *E. coli*). Určete, k jaké chybě došlo. Pro porovnání rozdílně predikovaných genů využijte program **EMBOSS Needle** (http://www.ebi.ac.uk/Tools/psa/emboss_needle/nucleotide.html).

SEKVENCE *E. COLI*

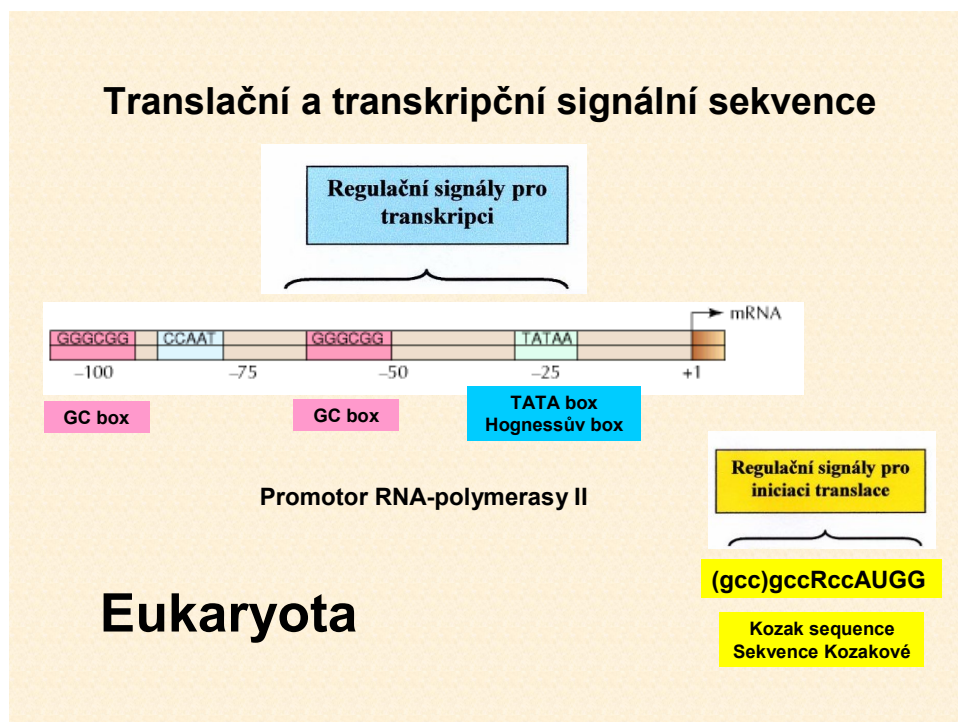
```
TTAAGAAATTCGGTATCAACTTCGAGGCCCTTTCAGGTCACCGTTGGCGTATTTCCGGGCGATAGATTTGTATTTGTAGCCGCAACTG
CGATCGCCCCCGGCGAAAAACGCAGATCCGGCGGTACTTTGTGAAATCACCGTTTCAATCCAGCCCAGCGTGCCGCGTGAACGAAACGAT
GGCGATCGTACAGTGTGCGGATCCAGACGTTCTGCGCCTGGAACCGGAGAAATCGACATCTGAGCCCCAGGCCGTGTTGGAGTAGTCGATAGA
GTAGCGTTGCGAGTCGCCCTGGTTGGCATCAGGCCACCACGAGAACCGGTGCGGCTAATCATACCCCAGGATAAAACAGCATCGTGGTGTG
GTAATTTACCCCTGAGTAAAGTGGTCGAGACTCCAGCGCAGGTTAATGGCAGCGTGCCAGCCGCTGGAGAGATCCAGTAGCGAGAAGCCACCA
CGGTAGTGGAGTCAGATTCGGTATCGTTTCAGGTCAGTGCCTTAAACCCGCCCTGCACCAAATAATATTGTTCCAGTGGATTCTTCAGCAGCGG
CATTTTTATAGCTGAAGTCGAGGGTCTGTTCCGGCGCGAAATACTGGTACTGGTGGTCAGACTGTGACCATAAGAGTTTCCACGGCTTTTTTC
CAGTAGCTTTACGCGCGGTCCACGTCCTGAGAGTAACCGACCCCGGTTTCGATGGTGTTCAGTTCGCGCGGAAACCACGCCCGTCAAGG
GTAATACTTTTCGTTTCGCGCGCTTTATCAAATTTGGAGCCACTACGACCGAGTTAAACCAGCCGGTAGCAGAAAGTCGACGGTTAAGTTCGCG
CAGATCTTTCGATTCGATCATCGCCCTTTTAAACGGCACCAGATTTGACAGTATTTCATCGCGGATTTGTGATCCTTCAAAGGTCACATGC
CCAAAGCGGTAACGTTCCGCACTGTTATAATCAATATCCAGAAAGGCTTTATGCAGGCCGAGCGCAATGCCAGCTGCGCTTTGGTAAATTCGC
TATCGAAATAACCTTTAGCAACGCAATGCTGGTTAAGGACTTTTTGAAATTTTCATAATCGCCCTGGTTCAGTACCGTGCCAATAGCCGGGCG
AGTATCGAGCAATTTCAAATAGTCTTTATCGGTCGCGCGCGCGCGCAATACCACATCGGTGCCGCCAATTAACACCGGCACGCCCTGGCGTG
ACTTTGGCGGATCAATACCTGCCGCCCTTTCTTTGGCGGTGGACGGAGATCAAATTCATGGTTCGGTGGTAATAACCCAGCGCTTCAGACCTT
CGCGGATGGCATCATCGACGCTGCGCGAAAGCGACGGTCTGGCGTCACTTTCATCACTTTCAATCGTAGAAAGTGCAGCAGAACGTTCTTTTC
CAGCTGTCCCATAACCCCTCGACTGTAGACGGACGTTTCGCGCGGACGGCAGATCCGTTAAGCAGAGTAAGCTTACACAGCATAACTGTCCG
ATATAGCGCATATTTCTCCTGAATATCCTTTTTCTCCTGCCCTGGAACCGCTTAAACAAAATCCAGTAATATGGATTAATAAAG
CAGACTAAACCCAAATATTTCTATGTTTTACTTTAGACCTATTCACGGTGGTATTGTTGTGCAAATACGCCCTTGTGTACAACCTTAACCCC
AATGACCGATTTTCGGGAGAGCGACACCATGAGTTTATTTGATAAAAAGCATCTGGTTTTCCCGCCGATGCCCTGGCAGTAAACACCCCG
ATGCCCGTAGCCAGCTGCATGCGGTCAACGGTCACTCAATGACCAATGTACCTGACGGAATGGAGATTGCCATTTTGGGATGGGTTGTTTCT
GGGGTGTGGAGCGTCTGTTCTGGCAGTTACCCGCGCTTTACAGCACCGCCGACGGCTATACCGCGGCTATACGCCAAATCCGACTTATCGGGA
AGTGTGCTCCGGTATGACCGGTATGCCGAAGCGGTACGCATTGTTTACGATCCTTCTGTATGAGCAGTATGAGCAGTATGATGATTTTGG
GAGAATACAGATCCCGCCAGGGAATGCGTCAGGGCAATGACCGCAATGACCGCATCGATTCGTTACGCGATTATCCGCTGACCCGAAACAGGATG
CCGACGCTCGCGCAGTCTGGAACGTTTTCAGGCGCGATGCTTCCGCCGATGACGATCGTCACATCACCCGAAATCGTAAACGCCACACC
GTTTTATATGCGGAAGATGACCACCAGCAATATCTGCATAAAAACCCGATATGGTTACTGTGGAATTGGCGGAATTGGCGTCTGTGCCACCG
GAAGCATAGCGTTACGGGTACAAATGTAGATTGTTGATAAAGTGCCTTTGTTTATGCCGAATGCAGCGTGAATGCCTTACCAGCCCTACAAA
TCGTCCAAATCAATATATGCGAGGACTGCGTAGGCCATGCGCATAGCGCATAGCGCATAGCGCATAGCGCATAGCGCATAGCGCATAGCGCAT
CTCTGGCGACTTTACAGTACCTTACGCTATAC TAGCCACTGAAAATGCCGGATCACTTTCTTCGAATCGGCTTTCAATGTGATTTACACAAA
TTAATCAACTTCCCTTCCGAGGATCTGGCCTGAAAGTTCGGATAAGATATGTTAAACAGTATTTAGTCATACTCTGCTTGATCGCTGTAAAGT
CGCTTCTCTCGATGCTGAGATCTCACTTGCCTCCTCAGCAAAATCAAATTAACCTGCTGGCTGATGAAGCAATATAAATGCCCAACCGG
TTCTGAATATGACGAAATAATCCGCGATGTTCTTTACGGTGGTCCAAATCGGTCGTAACCGCAGTTGCGATTCTCGCGGATCGTCCGCTGATG
GGCATTCTTCCAGCTTTTACAGCCTGTTCTCCGCTATATGTCGGCAGAACCTCTGAGCAACTGAGCTTTATCTCTCTTTCTCGTTAGTG
ACTGGCATGTTTATCCTGTTTGGGATTTAACCCGAAACGCATCGGTATGATTGCGCCAGAAGCGGTGGCTTTGCGTATCATCAACCCGATGC
GCTTTCGCTGTACGTTGCAACCCCGCTGGTGTGGTCTTCAACGGCTGGCGAACATGATCTCCGTTATTTCAAACCTGCCAATGGTACGTA
AGATGACATCACTCTGATGACATCTACGCGGTAGTGAAGCCGCTGCGCTGCGGCGGCTGTTACGTAACACAGGAACACGAGATGATGAAAC
GTCTTTGAGCTGGAATCCCGTACCGTTCGCTTCAATGACACCGCTGAAAACGTGATTTGGTTTGTATCTCCACGAGATGAGCAAAGCCTGA
AGAATAAGGTGGCGGAACATCCGCACTTAAGTTCCTCGTCTGTAATGAAGATATTGACCACATCATCGGTTATGTGATTTCAAAGACCTGCT
```


CACTATTATAACAATTCAATTGAAGAGCATTACAGTATAAAAATCCAATATCTAATAATCAAATGTGCCAATAATCGATTGATATCTCGTCAA
 GTGTTCAATTTGTTAGTAAATGATGATGATTTCACTTTGGTATTGAATAGAATTACCAATCACCTAGCATAAAATGAATATATTTGGAATTAT
 TAGAAGATATTAATGATATCCATTGATTATAAAACCCTTGGATTACCTGATTATCTAATACCAACTTGGTGGATCAATTGGGTAATAAGTT
 GGATTCTTTGGGTGATTTACAATTAGTATATTCATCTTCTACATCTAATAGGTTGAGAAATTTTATAATTGATATCCGAAAAACGATTTACCA
 AAATTTGCTAGTGATGGGAAATTTATATCAAGATATTGTTAAATTCATGTATAATAATACCAAAAATAAAATTTGAAAAATTAAGAAATAGAAAAAT
 TTATTAATAATTAATTAATATTGCTAGTGATGGGAAATTTAAATTAATAATTAATGATAATCAATTAATTTGGTTTTTATTGATTGAATCTATAAA
 ATCATCTTTACAATAAGTAGTAAGTGTACGTGGTCAATAAAGTAATGATATATATGTGCGTGTACTGAGTGATGATTTTCTATTTTATCTCGCGC
 GTTGTATTTTTTGTGTTATTGTT
 ATTTTTTTTTTCTCTCTCTCTCCCACTTTCTACTCTTTGAACATGTGGTGAAAAAATAAATAGAACAGATCTTTTTTATCCATATAT
 TTCAAACCGACTTGTTCCTTTTTTTCTTTTACAACCAATTAATACAAAAGAAAGACAACAATCATAATGTCTACTTCTGTTGAACCCAATGAA
 ACAGAAGCTTTGTTGAGAAAGCAGAATGATCTTTCCACAACCTGCTCAATTGAAGAAAAATATCCTCACCAACAGGAGAGGCTGCAGAAGACG
 ATGACGACACTCTTAAAAGAACCAATATGATGAAGCTAAAGAAACCGCTGAATCTTTAAAACAAGTTGAATCGATATTAGCACCTATTGTTTT
 CACTGCATGTCAATTTTCGTGAGATTTTATCGTATTTTCAAGTAAATGACCATGTTGTTGGGATGAAGCTCATTTTGGTAAATTTGGATCCTAT
 TATTTACGACACGAATTTTATCATGATGTTTCACTCTCCATTGGGTAATGTTAGTTGGTTTATCTGGTTATTTGGCTGGATACAATGGATCTT
 GGATTTCCCAAGTGGTAAAAATACCTGATTATATTGATTACTAAAAATGAGATTGTTAATGCCACTTTCTGCTGCTTGTGTGACCATT
 GGCTTATTTCACTGGGAAAGAAGTTGGATTTTCCATGTTTACTACTTGGTTATTTACTTTGATGGTGGCTCTTGAATCAAGTTATGTCACTTTA
 GGTAATTCATTTGTTGGATTCAATGTTGTTATTTCTCACCGTTGCTACTGTTTCTGTTTTTCACGTTTCAACAATTTTAAACAATAAATCAC
 AAGAATTTCTAGAAAATGGTGGAAATGGATCTTTTAACTGGTGTTCATTGGTTGTTGTTGTTGTTGTTGTTGTTGTTGTTGTTGTTGTTGTTGTT
 CACATTTGGTCCGGATTTTCACTGTTGTTGACCTTTGGAATAAATGAGTGATAAATCTATTTTCATGGACAAAATACATTCAACATTTGGTTGCT
 AGAATTTGTTGCTTTGATTCTTGTCCCAATTTTCAATTTTCAATGCTTTTCAATTTAAAGTTCAATTTTGAATTTGTTGATAAATCGGGTACTGGTGATG
 CCAATATGTCATCACTTTTCCAAGTAATTTGGCTGGTTCCGATGTTGGTGGTGGCCACGTAAGTATCCATGTTCCACTCGGTTATCACTTT
 AAAGAATCAAGGTTTAAAGTGGTGGCCTTTTACTCTCCACGTTCAACATTTCCAGAAGTTCAAACAACAACAAGTTACTACTTATGGTCCAC
 AAAGATTCAAAACAACAATTTGATTTTCCAAAGCTAGAGGACAACCTTATTATGATACTTCTGGTAACACCACGACATGAATATATTTTTG
 ACGGTATGCATGTAAGATTGATGCATCCACAACTGGTAGAACTTACATACTCATGATATCCAGCTCCAGTGTCTAAATCTGAATATGAAGT
 TGCATGTTATGGTAATTTGACTATTGGTGATCCTAAAGATAATTTGGACTGTTGAAATTTGGAACAAGCAAGTATGAAGATAAATAGATTA
 CATCTTTGACTTCGTCAATTTAGATTGAAGAATGAAGTGAATGTTATTTGGGGTCACTGGTACTACATTACCTCAATGGGGTTCAGAC
 AAGGTGAAGTTGTTGTTTACAAGAACCATTAAAAAGACAAGAGAACCTT

ÚKOL 3

Charakterizujte část genomu *Candida albicans* (kvasinka) rovněž pomocí heuristického modelu.

PREDIKCE GENŮ U MNOHOBUNĚČNÝCH EUKARYOT



Mnohobuněčná eukaryota se vyznačují komplexní organizací genomu, geny jsou separovány dlouhými intergenovými úseky, geny obsahují mnoho intronů, i velmi dlouhých. Exony/introny jsou identifikovány pomocí míst setřihu (GT na 5'konci intronu, AG na 3'konci). Vzniká velké množství chyb! Dlouhé introny jsou určeny jako intergenové úseky, krátké intergenové úseky jako introny, krátké exony nemusí být identifikovány.

ÚKOL 5

Analyzujte část genomu mnohobuněčného eukaryotického organismu také pomocí programu **GENSCAN** (<http://genes.mit.edu/GENSCAN.html>). Porovnejte s výsledky z úkolu 4.