### Hidden Markov models in gene finding

Broňa Brejová Bioinformatics research group David R. Cheriton School of Computer Science University of Waterloo

# **Topics for today**

- What is gene finding (biological background)
- HMMs in gene finding tools of the trade:
  - Viterbi algorithm
  - Higher order states
  - State transition diagram

### What is a genome?

cggtgaaactgcacgattgttgctggcttaaagatagaccaatcagagtgtgtaacgtca tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca tgggcgtatttgcgctagtgttgggtgttccgctgtgctgtttttccgtcatggctcgca  ${\tt ctaagcaaactgctcggaagtctactggtggcaaggcgccacgcaaacagttggccacta}$ aggcagcccgcaaaagcgctccggccaccggcggcgtgaaaaagccccaccgctaccggc cgggcaccgtggctctgcgcgagatccgccgttatcagaagtccactgaactgcttattcgtaaactacctttccagcgcctggtgcgcgagattgcgcaggactttaaaacagacctgc gtttccagagctccgctgtgatggctctgcaggaggcgtgcgaggcctacttggtagggc tatttgaggacactaacctgtgcgccatccacgccaagcgcgtcactatcatgcccaagg acatccagctcgccgccgcatccgcggagagaggggggtgattactgtggtctctctgac ggtccaagcaaaggctcttttcagagccaccaccttttcaagtaaagtagctgtaagaaa ccaatttaagacaaaagggaatgcattgggagcacttttcgttttaatgctactgaaggc

Human genome: 3 billion letters (sequenced in 2001) What do they mean?

# What are genes?



Genes are portions of genome encoding proteins

# **Closer look on genes**

### **Process of protein production:**



**Translation:** triple of letters in mRNA  $\rightarrow$  one amino acid in protein



### Task of gene finding

cggtgaaactgcacgattgttgctggcttaaagatagaccaatcagagtgtgtaacgtca tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca tgggcgtatttgcgctagtgttgggtgttccgctgtgctgtttttccgtcatggctcgca ctaagcaaactgctcggaagtctactggtggcaaggcgccacgcaaacagttggccacta aggcagcccgcaaaagcgctccggccaccggcggcgtgaaaaagccccaccgctaccggc cgggcaccgtggctctgcgcgagatccgccgttatcagaagtccactgaactgcttattc gtaaactacctttccagcgcctggtgcgagattgcgcaggactttaaaacagacctgc gtttccagagctccgctgtgatggctctgcaggaggcgtgcaggcgtgaggcctactggtagggc tatttgaggacactaacctgtgcgcgagagagggcgtgatactgtgtgccaagg acatccagctcgccgcatccgcggagagggggggggtgattactgtggtctctgcgagg acatccagctcgccgcatccgcggagagggggggggtgattactgtggtctctctgac

# Task of gene finding



# Toy Hidden Markov model for gene finding





p(q,y): joint prob. over genomic sequences y and gene structures q



# Toy Hidden Markov model for gene finding







# HMMs for gene finding: training





### Training:

- Training set of sequences with known gene structures
- Fully observed model (both q and y known)
- Maximum likelihood estimation by relative frequencies

# HMMs for gene finding: inference





### Inference: Viterbi algorithm

- Genomic sequence y observed, state sequence q unknown
- Viterbi algorithm computes state sequence q\*

with highest probability given y:

$$q^* = \arg \max_q p(q|y) = \arg \max_q p(q,y)$$

# HMMs for gene finding: higher order states

### Order 0:



# $\eta$ specifies $p(y_t|q_t)$

# Order 1:



# $\eta$ specifies $p(y_t|q_t,y_{t-1})$

.

# HMMs for gene finding: higher order states

 $\eta$  specifies  $p(y_t|q_t, y_{t-1})$ :

. . .

qt	Yt-1	a	С	g	t
	a	0.24	0.23	0.34	0.19
	С	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
	а	0.30	0.18	0.27	0.25
	С	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

# HMMs for gene finding: state transition diagrams







Graph representation of A:



# More complex HMMs: change state transition diagrams

Recall: triple of letters in exon  $\rightarrow$  one amino acid



# New states have separate emission probability





0.27 0.23 0.23

0.27

η	a	С	g	t	η	a	С	g	t
	0.24	0.27	0.28	0.21	0	0.26	0.26	0.32	0.16
	0.26	0.22	0.22	0.30	1	0.30	0.24	0.20	0.26
	0.27	0.23	0.23	0.27	2	0.17	0.32	0.31	0.20
						0.26	0.22	0.22	0.30

# More complex HMMs: change state transition diagrams

Keep consistent triples across introns:



### More complex HMMs: remember those signals?

Exor

#### Intron

ccatcccctatatttatggcaggtgaggaaagggtggggggctgggg attcatcatcatgggtgcatcggtgagtatctcccaggccccaatc agaagatctaccccaccatctggtaagtgtgtcccaccactgcccc acagagtgagcccttcttcaaggtgggtgtcagggcctccccc tgcagaacctcatggtgctgaggggccaagcctgggccgggg tcgatgaatttgggatcatccggtgagagctcttcctctctcgg agatgacgtccgtgatgagagggggggggggggcgcaccccagtcccca tttcttgtggctattttaaaaggtaattcatggagaaatagaaaaa tttcttgtggctattttaa**aaggtaatt**catggagaaatagaaaaa tttgaagaaactccacgaagaggttgatggcagtgactttcggaaa agtggatgcccttaaaggaaccgtggagtaccaaccccctgcagt cgacccgtgaccctcgtgagggggggggctc aaatgcagtggaagagggactagtacgtgagccatgctgggagtgt catggcgggtgtgctgaagaaggtgagacgaatggaggtcactgtt gaagctggctatggttaaagcggtaagtagctaagtcagttttgtt attagaagaggtgattcttcaggtaaagaaaagttgactatttag

### More complex HMMs: remember those signals?

ExonIntronccatcccctatatttatggcaggtgaggaaagggtggggggctggggattcatcatgggtgcatcggtgagtatctcccaggccccaatcagaagatctacccccaccatctggtaagtgtgtgcccccactgccccacagagtgagcccttcttcaaggtgggtggtgtgtcagggcccccca

Insert a series of states between exon and intron:



# **Example of an HMM for gene finding**



# Summary

- HMMs are useful tool for gene finding
- Viterbi algorithm used for inference
- Higher order states model dependencies among adjacent letters
- Complex state transition diagrams incorporate biological knowledge

# This is just a beginning

- More complex models used in practice: lengths of exons/introns, complex signal models, experimental information, multiple species,...
- HMMs for other bioinformatics tasks: protein secondary structure, protein families, segmenting genome,...