Advanced Algorithms and Models for Computational Biology

-- a machine learning approach

Computational Genomics II:

HMM variants and Comparative
Gene Finding



Eric Xing

Lecture 5, February 1, 2005

Reading: Chap 3, 5 DEKM book Chap 9, DTW book

Higher-order HMMs



- The Genetic Code
 - 3 nucleotides make 1 amino acid
 - Statistical dependencies in triplets
- Question:
 - Recognize protein-coding segments with an HMM

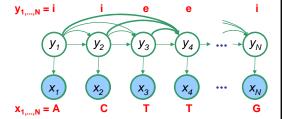
		U	С	Α	G
	U	UUU phe UUC Phe UUA 1eu	UCU UCC UCA ser UCG	UAU UAC UAA Stop UAG Stop	UGU cys UGC cys UGA Stop UGG Stop
	0	CUU CUC CUA CUG	CCU CCC CCA pro CCG	CAU his CAC CAA gln CAG	CGU CGC arg CGA CGG
	A	AUU ile AUA TAUG met	ACU ACC ACA ACG	AAU asn AAC AAA AAG lys	AGU ser AGC AGA arg AGG
	G	GUU GUC GUA GUG	GCU GCC _{ala} GCA GCG	GAU asp GAC GAA GAA glu	GGU GGC GGA GGG

Higher-order HMMs



- Every state of the HMM emits 1 nucleotide
- Transition probabilities:

Probability of a state at one position, given those of 3 previous positions (triplets): $P(y_i | y_{i-1}, y_{i-2}, y_{i-3})$



- Emission probabilities: $P(x_i | y_i)$
- Algorithms extend with small modifications

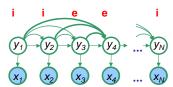
Inference on Higher-order HMMs

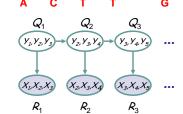


- Building 1st-order HMM on "mega" state
- Use FB algorithm as usual



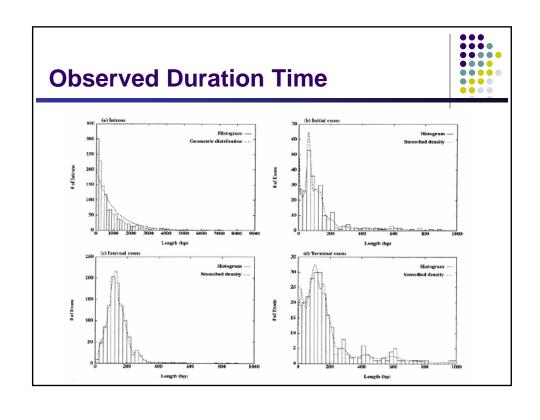
- $P(Q_2|R)$
- $\rightarrow P(Y_2, Y_3, Y_4 | X)$





 $\rightarrow P(Y_3|X)=\Sigma_{y_2,y_4}P(Y_2, Y_3, Y_4|X)$

• Length distribution of region X: E[I_X] = 1/(1-p) • Geometric distribution, with mean 1/(1-p) • (homework: derive this) • This is a significant disadvantage of HMMs • Several solutions exist for modeling different length distributions



Poisson Point Process



- A counting process that represents the total number of occurrences of discrete events during a temporal/spatial interval
 - the number of occurrences in any internal of length τ is Poisson distributed with parameter $\lambda \tau$:

$$p(A(t+\tau)-A(n)=n)=e^{-\lambda\tau}\frac{(\lambda\tau)^n}{n!}$$

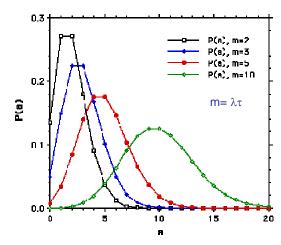


- the number of occurrences in disjoint intervals are independent
- the duration of the interval between two consecutive occurrences has the following distribution:

$$p(\tau < s) = 1 - e^{-\lambda s}$$

Poisson point process





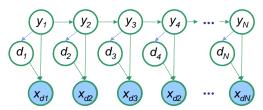
Truncation is needed at both ends!

Generalized HMM



Upon entering a state:

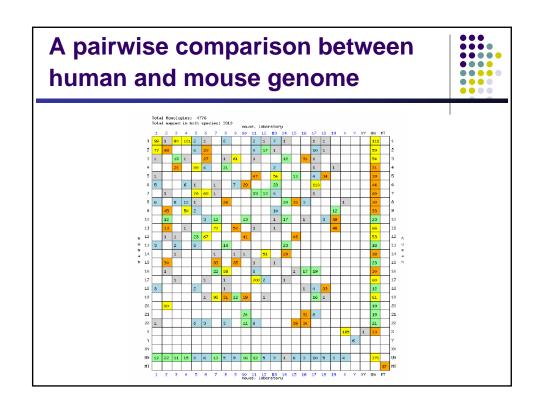
- 1. Choose duration d, according to probability distribution
- 2. Generate d letters according to emission probs
- 3. Take a transition to next state according to transition probs

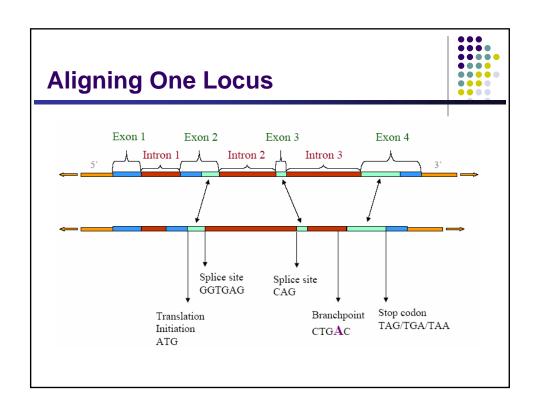


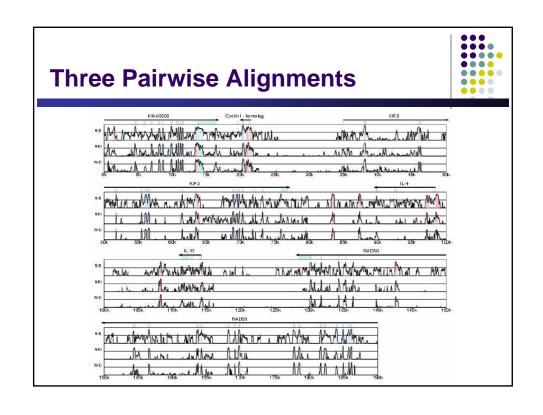
Disadvantage: Increase in complexity:

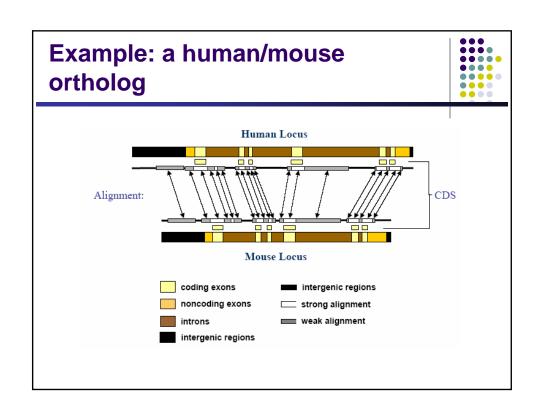
Time: O(D²) Space: O(D)

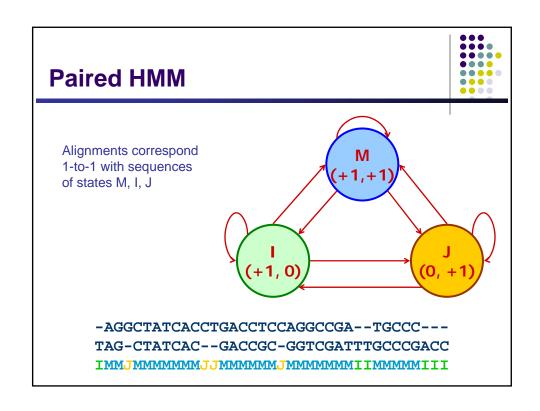
where D = maximum duration of state

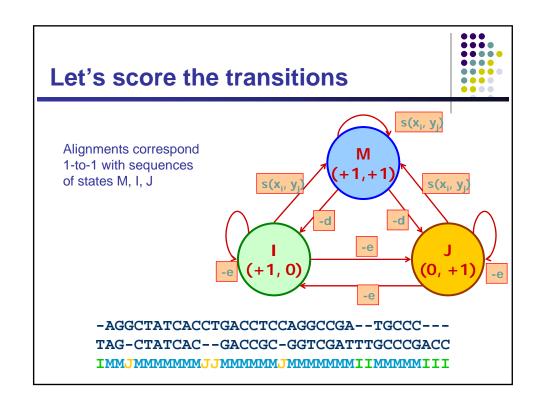


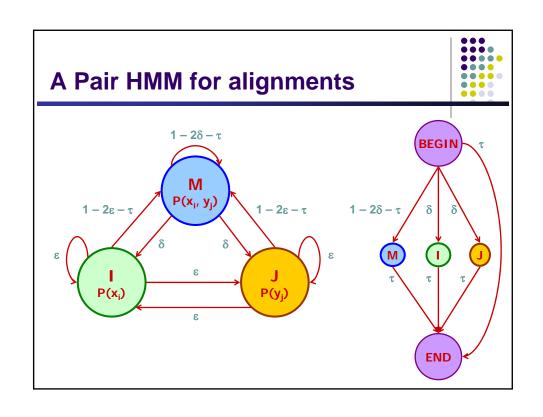


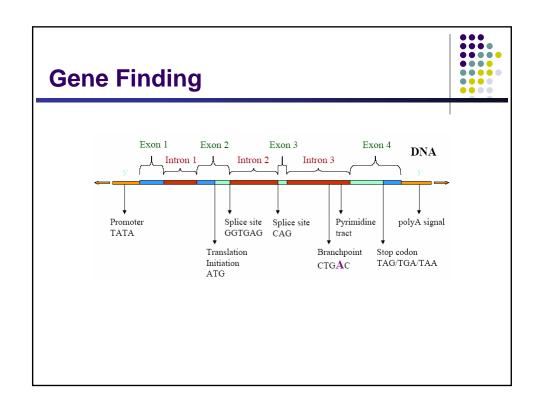


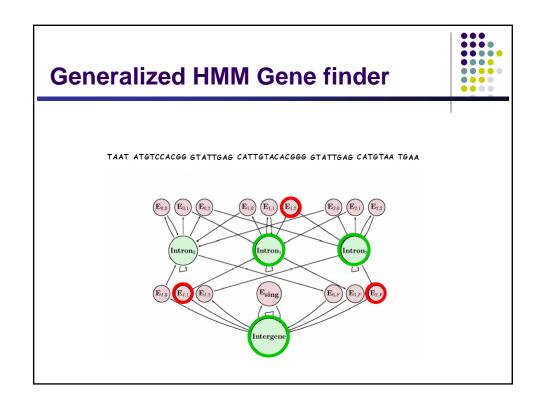


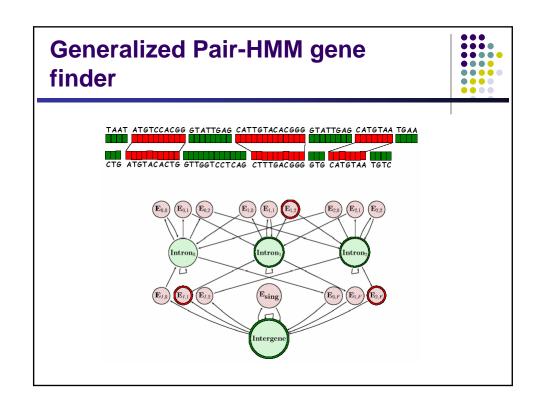


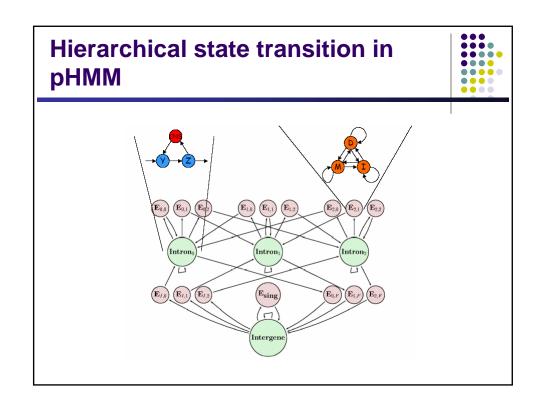


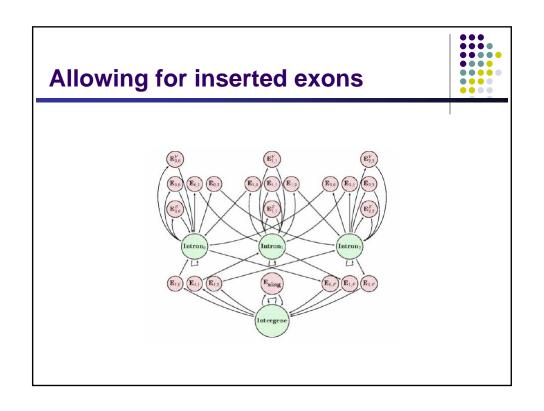












Acknowledgments



- **Serafim Batzoglou**: for some of the slides adapted or modified from his lecture slides at Stanford University
- Lior Pachter': for some of the slides modified from his lectures at UC Berkeley