

HMM variants

--- for comparative gene finding,
etc.



10-810, CMB lecture 4---Eric Xing

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	C	A	G
U	UUU phe UUC UUA leu UUG	UCU UCC ser UCA UCG	UAU tyr UAC UAA Stop UAG Stop	UGU cys UGC UGA Stop UGG Stop
C	CUU CUC leu CUA CUG	CCU CCC pro CCA CCG	CAU his CAC CAA gln CAG	CGU CGC arg CGA CGG
A	AUU AUC ile AUA AUG met	ACU ACC thr ACA ACG	AAU asn AAC AAA lys AAG	AGU ser AGC AGA arg AGG
G	GUU GUC val GUA GUG	GCU GCC ala GCA GCG	GAU asp GAC GAA glu GAG	GGU GGC gly GGA GGG

Higher-order HMMs



Every state of the HMM emits 1 nucleotide

Transition probabilities:

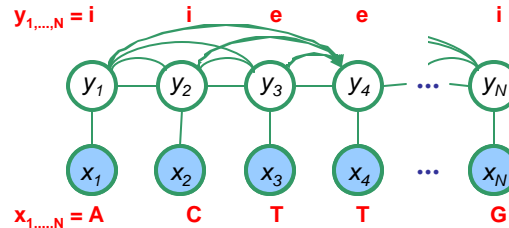
Probability of a state at one position, given those of 3 previous positions (triplets):

$$P(\pi_i | \pi_{i-1}, \pi_{i-2}, \pi_{i-3})$$

Emission probabilities:

$$P(x_i | \pi_i)$$

Algorithms extend with small modifications



Modeling the Duration of States

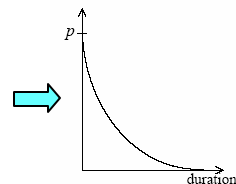


Length distribution of region X:

$$E[l_X] = 1/(1-p)$$

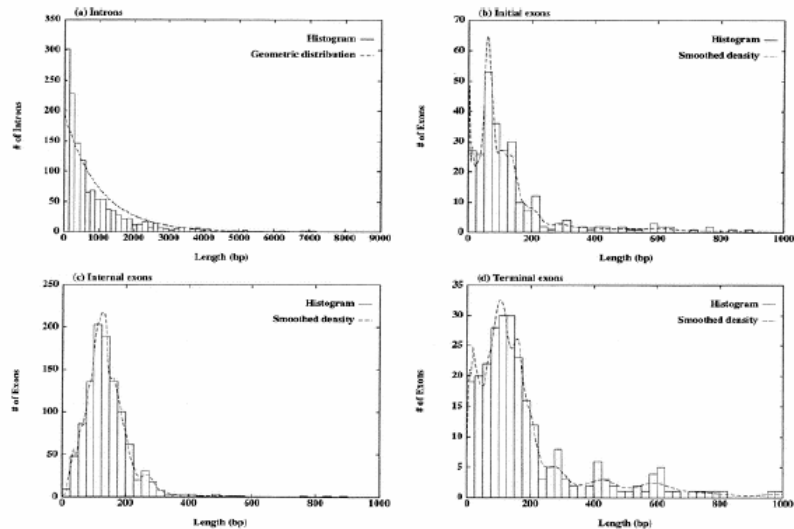
- Geometric distribution, with mean $1/(1-p)$

This is a significant disadvantage of HMMs



Several solutions exist for modeling different length distributions

Observed duration time



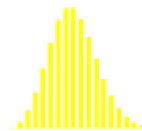
Poisson point process



A counting process that represents the total number of occurrences of discrete events during a temporal/spatial interval

- the number of occurrences in any interval of length τ is **Poisson distributed** with parameter $\lambda\tau$:

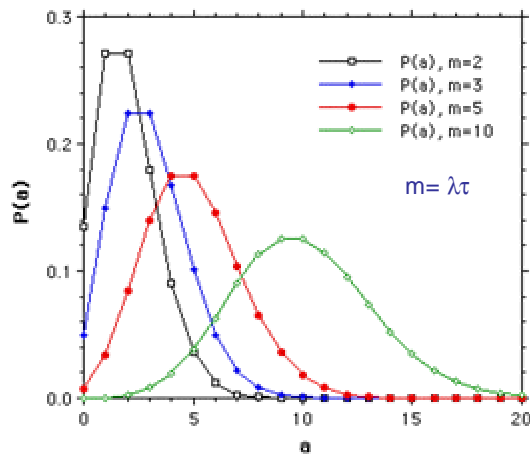
$$p(A(t + \tau) - A(t) = n) = e^{-\lambda\tau} \frac{(\lambda\tau)^n}{n!}$$



- 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



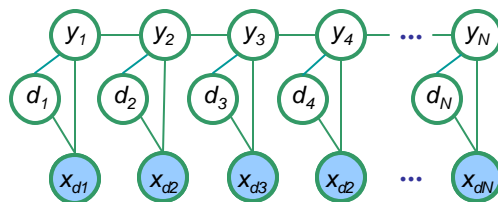
Truncation is needed at both ends!

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

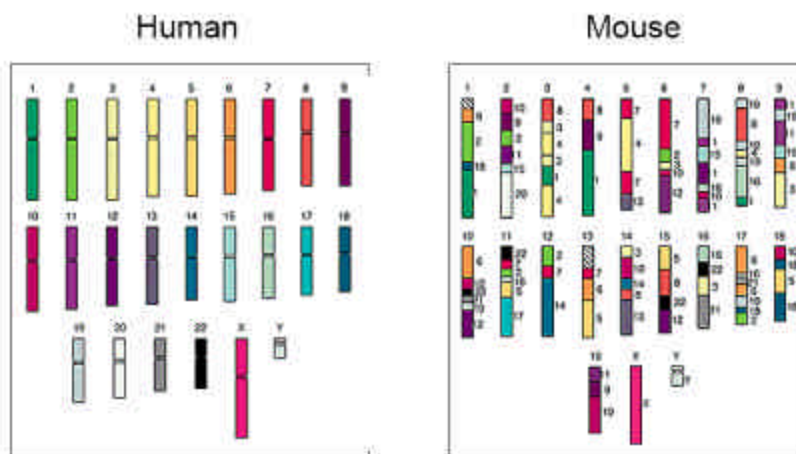


Disadvantage: Increase in complexity:

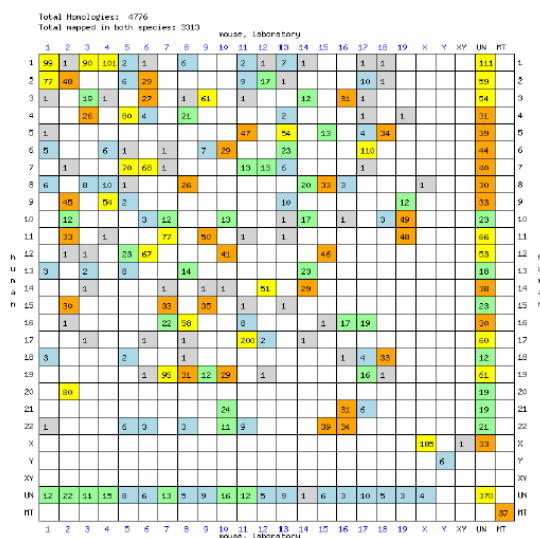
Time: $O(D^2)$
Space: $O(D)$

where D = maximum duration of state

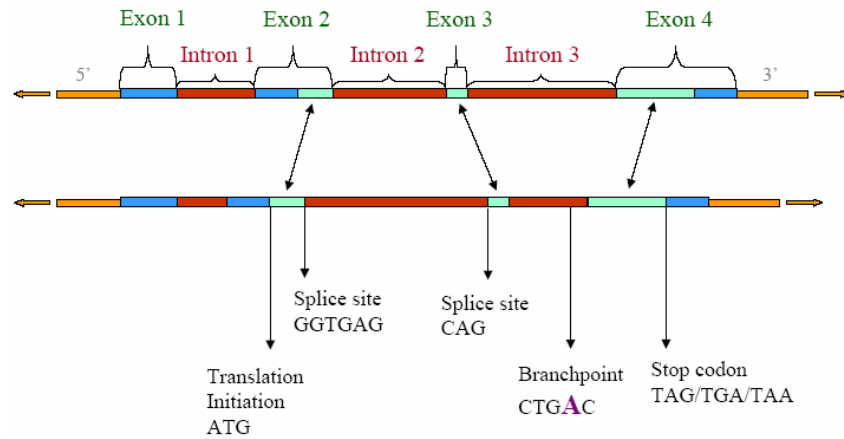
Comparative genomics



A pairwise comparison between human and mouse genome



Aligning one locus



Pairwise alignment - a close-up view



```

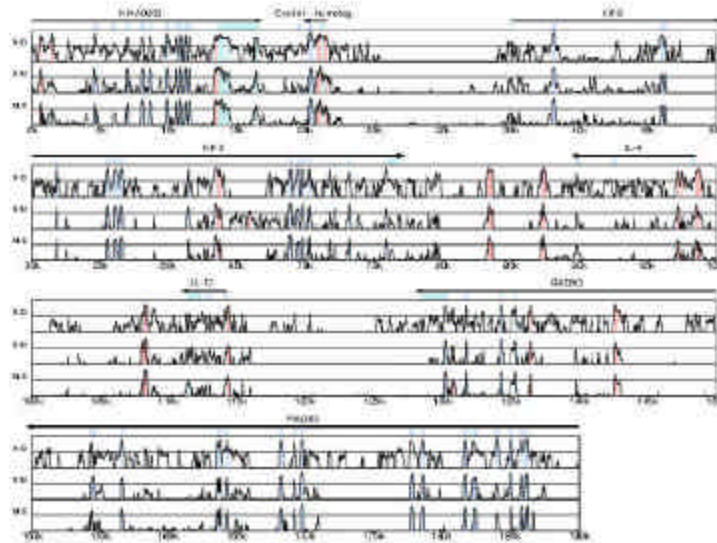
50      .   :   .   :   .   :   .   :
247 GGTGAGGTCGAGGACCCTGCA CGGAGCTGTATGGAGGGCA AGAGC
   | :  ||  ||| :  |||  --:|  ||  |:|  |||---|||
368 GAGTCGGGGGAGGGGGCTGCTGTTGGCTCTGGACAGCTTGCATTGAGAGG

100     .   :   .   :   .   :   .   :
292 TTC          CTACAGAAAAGTCCCAGCAAGGAGCCACACTTCACTG
   ||-----||  |  |:|  |:  |||::|:-||  ||:|  |
418 TTCTGGCTACGCTCTCCCTTAGGGACTGAGCAGAGGGCT CAGGTCGCGG

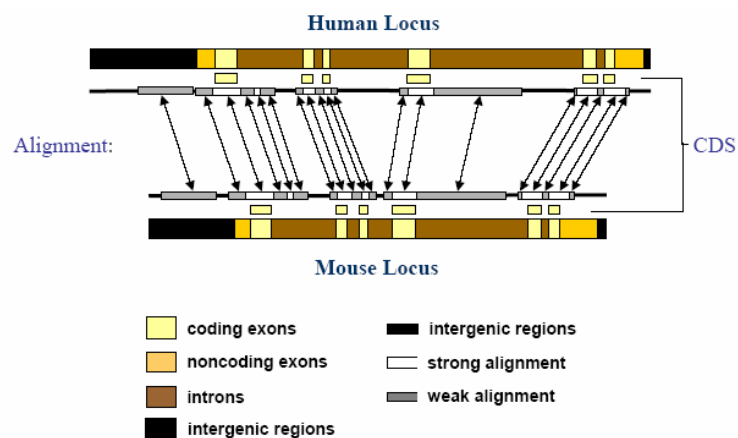
150     .   :   .   :   .   :   .   :
332          ATGTCGAGGGGAAGACATCATTGGGATGTCAGTG
   -----|||:|||:|||:|||:|||:|||:|||
467 TGGGAGATGAGGCCAATGTCGAGGGGAAGACATCATTGGGATGTCAGTG

200     .   :   .   :   .   :   .   :
367 TTCAACCTCAGCAATGCCATCATGGGCAGCGGCATCCTGGGACTCGCCTA
   |||:|||:|||:|||:|||:|||:|||:|||:|||:|||:|||
517 TTCAATCTCAGCAACGCCATCATGGGCAGTGAATTCTGGGGCTCGCCTA
    
```

Three pairwise alignments



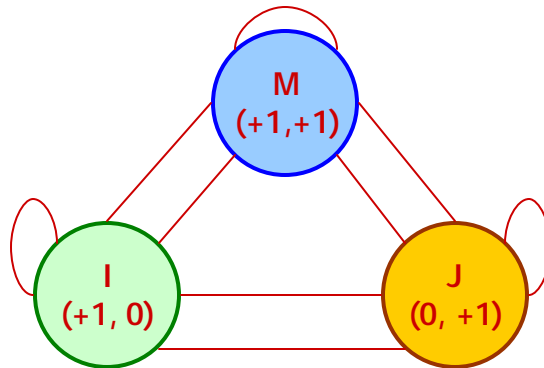
Example: a human/mouse ortholog



Paired HMM



Alignments correspond
1-to-1 with sequences
of states M, I, J

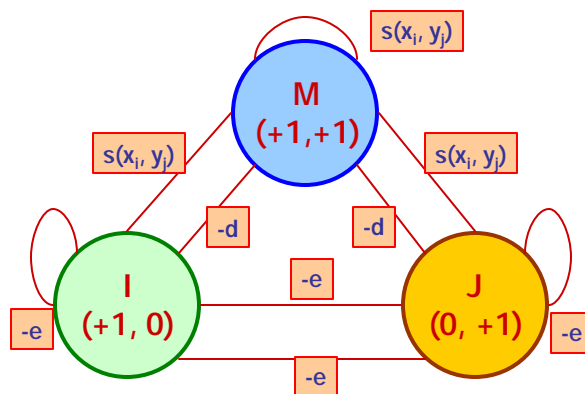


-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC--
TAG-CTATCAC--GACCGC-GGTCGATTGCCCCGACC
I M M J M M M M M M M J J M M M M M M J M M M M M M M I I M M M M M I I I

Let's score the transitions

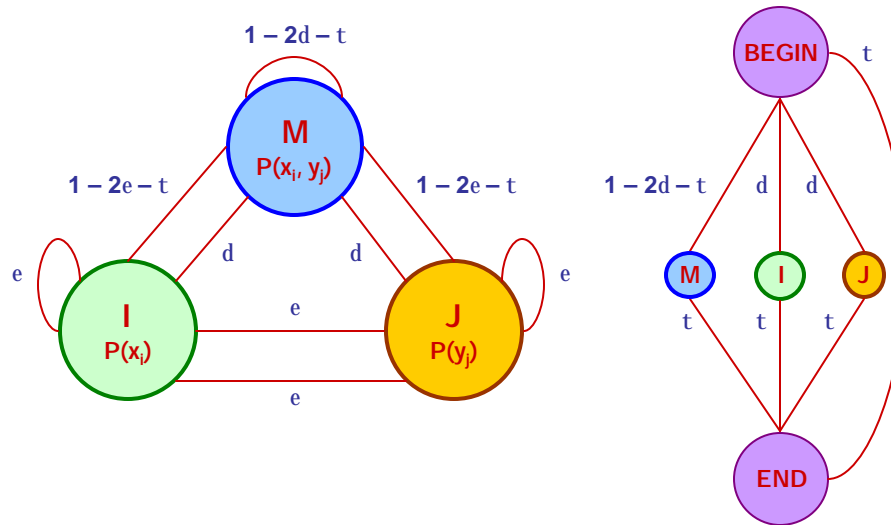


Alignments correspond
1-to-1 with sequences
of states M, I, J

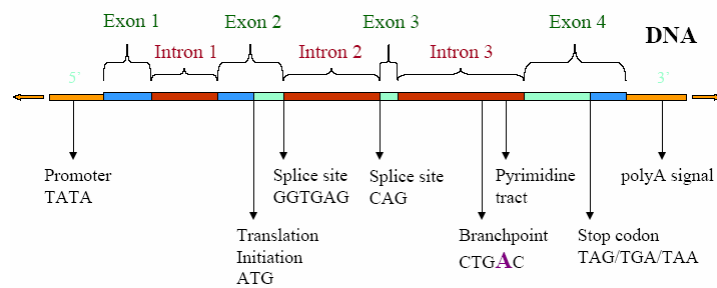


-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC--
TAG-CTATCAC--GACCGC-GGTCGATTGCCCCGACC
I M M J M M M M M M M J J M M M M M M J M M M M M M M I I M M M M M I I I

A Pair HMM for alignments



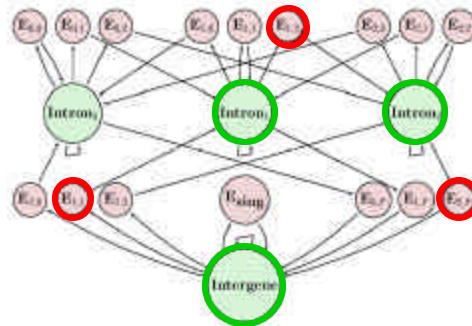
Gene finding



Generalized HMM Gene finder



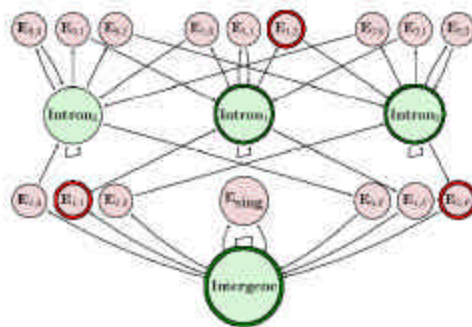
TAAT ATGTCCACGG GTATTGAG CATTGTACACGGG GTATTGAG CATGTAA TGAA



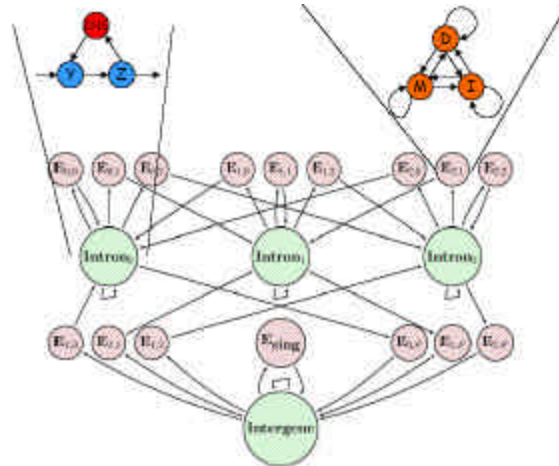
Generalized Pair-HMM gene finder



TAAT ATGTCCACGG GTATTGAG CATTGTACACGGG GTATTGAG CATGTAA TGAA
 CTG ATGTACACTG GTTGGTCCTCAG CTTTACACGGG GTG CATGTAA TGTC



Hierarchical state transition in pHMM



Allowing for inserted exons

