

# CSE P 590 A Computational Biology

## Gene Prediction

## Gene Finding: Motivation

Sequence data flooding into Genbank

What does it mean?

protein genes, RNA genes, mitochondria,  
chloroplast, regulation, replication, structure,  
repeats, transposons, unknown stuff, ...

More generally, how do you learn from  
complex data in an unknown language

8

## Protein Coding Nuclear DNA

Focus of this lecture

Goal: Automated annotation of new seq data

State of the Art:

In Eukaryotes:

predictions ~ 60% similar to real proteins  
~80% if database similarity used

Prokaryotes

better, but still imperfect

Lab verification still needed, still expensive

Largely done for Human; unlikely for most others

9

## Biological Basics

Central Dogma:

DNA [transcription](#) RNA [translation](#) Protein

Codons: 3 bases code one amino acid

Start codon

Stop codons

3', 5' Untranslated Regions (UTR's)

10

# RNA Transcription

(This gene is heavily transcribed, but many are not.)

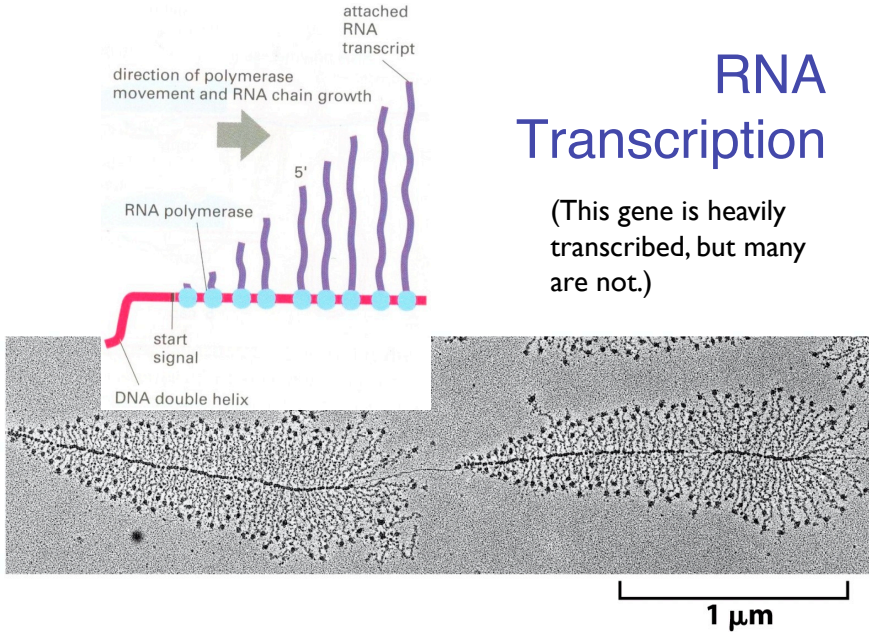
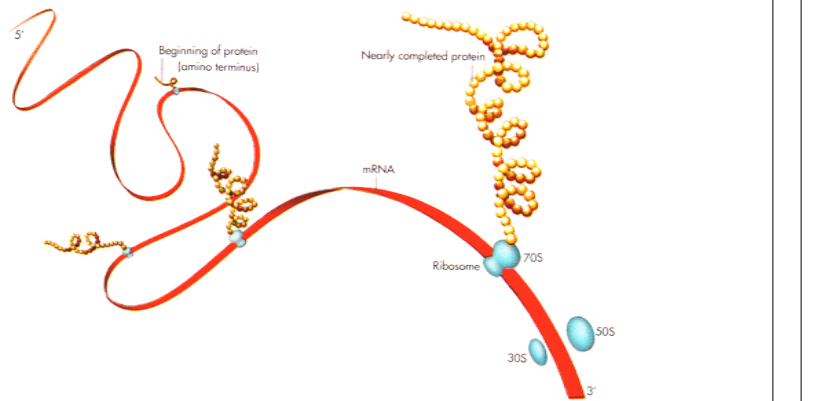


Figure 6-9 Molecular Biology of the Cell 5/e (© Garland Science 2008)

# Codons & The Genetic Code

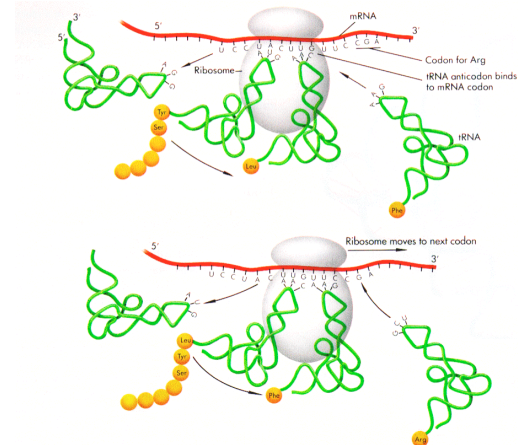
		Second Base				Third Base	
		U	C	A	G		
First Base	U	Phe	Ser	Tyr	Cys	U	Ala : Alanine
		Phe	Ser	Tyr	Cys	C	Arg : Arginine
		Leu	Ser	Stop	Stop	A	Asn : Asparagine
		Leu	Ser	Stop	Trp	G	Asp : Aspartic acid
C	Leu	Pro	His	Arg	U	Cys : Cysteine	
		Pro	His	Arg	C	Gln : Glutamine	
		Pro	Gln	Arg	A	Glu : Glutamic acid	
A	Ile	Thr	Asn	Ser	U	Gly : Glycine	
		Thr	Asn	Ser	C	His : Histidine	
		Thr	Lys	Arg	A	Ile : Isoleucine	
		Met/Start	Lys	Arg	G	Leu : Leucine	
G	Val	Ala	Asp	Gly	U	Lys : Lysine	
		Ala	Asp	Gly	C	Met : Methionine	
		Ala	Glu	Gly	A	Phe : Phenylalanine	
		Ala	Glu	Gly	G	Pro : Proline	
						Ser : Serine	
						Thr : Threonine	
						Trp : Tryptophane	
						Tyr : Tyrosine	
						Val : Valine	

# Translation: mRNA → Protein



Watson, Gilman, Wilkowsky, & Zoller, 1992

# Ribosomes



Watson, Gilman, Wilkowsky, & Zoller, 1992

## Idea #1: Find Long ORF's

**Reading frame:** which of the 3 possible sequences of triples does the ribosome read?

**Open Reading Frame:** No stop codons

In random DNA

average ORF =  $64/3 = 21$  triplets

300bp ORF once per 36kbp per strand

But average protein ~ 1000bp

15

## A Simple ORF finder

start at left end

scan triplet-by-non-overlapping triplet for AUG

then continue scan for STOP

repeat until right end

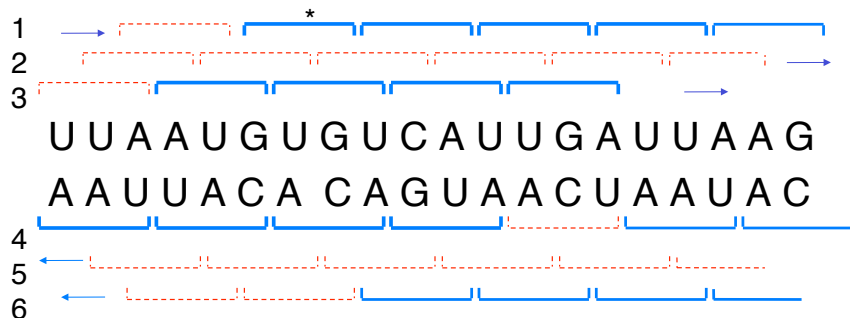
repeat all starting at offset 1

repeat all starting at offset 2

then do it again on the other strand

16

## Scanning for ORFs



17

## Idea #2: Codon Frequency

In random DNA

Leucine : Alanine : Tryptophan = 6 : 4 : 1

But in real protein, ratios ~ 6.9 : 6.5 : 1

So, coding DNA is not random

Even more: synonym usage is biased (in a species dependant way)

examples known with 90% AT 3<sup>rd</sup> base

Why? E.g. efficiency, histone, enhancer, splice interactions

18

## Recognizing Codon Bias

Assume

Codon usage i.i.d.;  $abc$  with freq.  $f(abc)$

$a_1a_2a_3a_4 \dots a_{3n+2}$  is coding, unknown frame

Calculate

$$p_1 = f(a_1a_2a_3)f(a_4a_5a_6) \dots f(a_{3n-2}a_{3n-1}a_{3n})$$

$$p_2 = f(a_2a_3a_4)f(a_5a_6a_7) \dots f(a_{3n-1}a_{3n}a_{3n+1})$$

$$p_3 = f(a_3a_4a_5)f(a_6a_7a_8) \dots f(a_{3n}a_{3n+1}a_{3n+2})$$

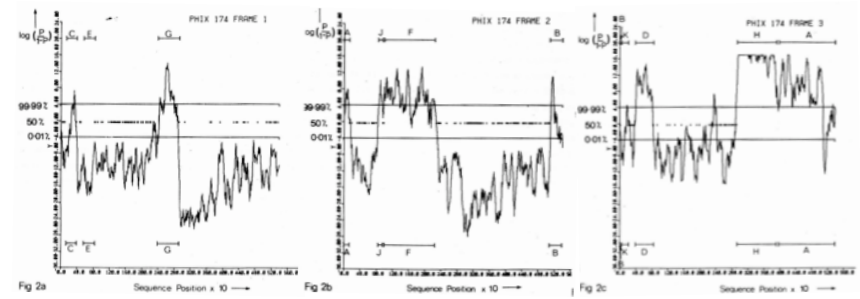
$$P_i = p_i / (p_1 + p_2 + p_3)$$

More generally:  $k$ -th order Markov model

$k = 5$  or  $6$  is typical

19

## Codon Usage in $\Phi$ x174



Staden & McLachlan, NAR 10, 1 1982, 141-156

22

## Promoters, etc.

In prokaryotes, most DNA coding

E.g.  $\sim 70\%$  in *H. influenzae*

Long ORFs + codon stats do well

But obviously won't be perfect

short genes

5' & 3' UTR's

Can improve by modeling promoters, etc.

e.g. via WMM or higher-order Markov models

23

## Eukaryotes

As in prokaryotes (but maybe more variable)

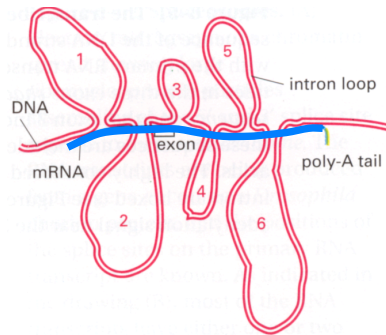
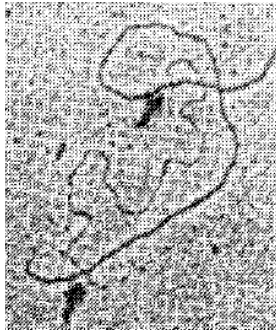
promoters

start/stop transcription

start/stop translation

24

And then...



Nobel Prize of the week: P. Sharp, 1993, Splicing

## Mechanical Devices of the Spliceosome: Motors, Clocks, Springs, and Things

Jonathan P. Staley and Christine Guthrie

CELL Volume 92, Issue 3 , 6 February 1998, Pages 315-326

Figure 2. Spliceosome Assembly, Rearrangement, and Disassembly Requires ATP, Numerous DExD/H box Proteins, and Prp24. The snRNPs are depicted as circles. The pathway for *S. cerevisiae* is shown.

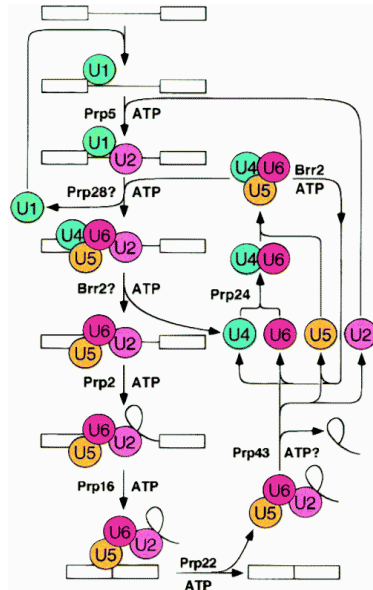
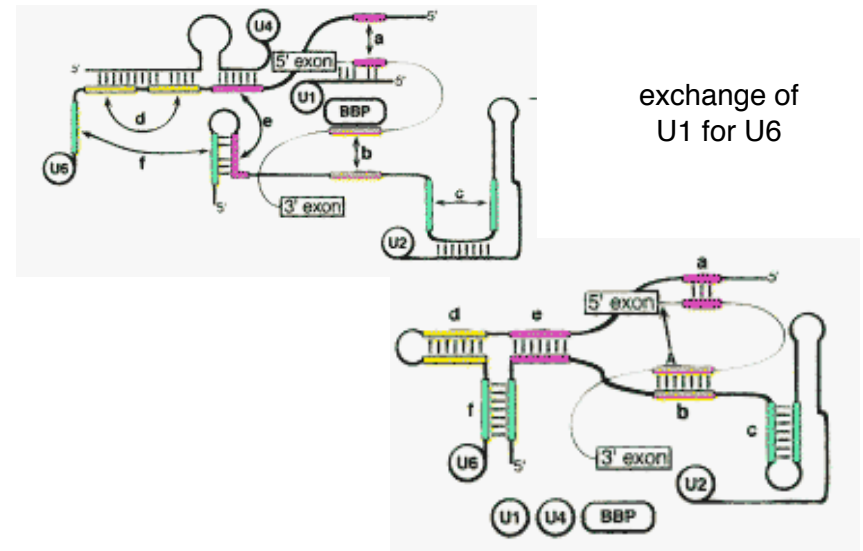


Figure 3. Splicing Requires Numerous Rearrangements



exchange of U1 for U6

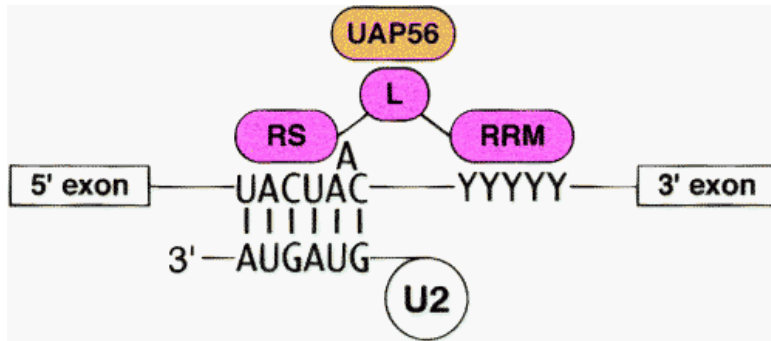
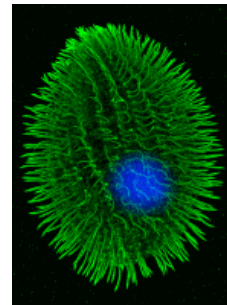


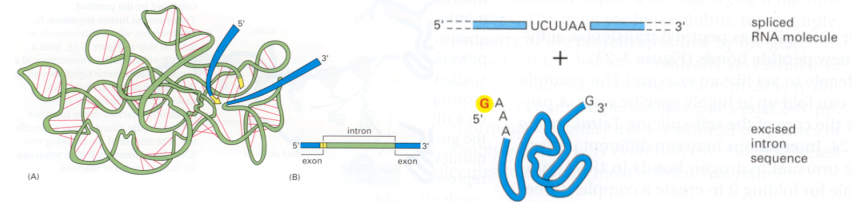
Figure 6. A Paradigm for Unwindase Specificity and Timing? The DExD/H box protein UAP56 (orange) binds U2AF65 (pink) through its linker region (L). U2 binds the branch point. Y's indicate the polypyrimidine stretch; RS, RRM as in Figure 5A. Sequences are from mammals.

33

## Hints to Origins?



Tetrahymena thermophila



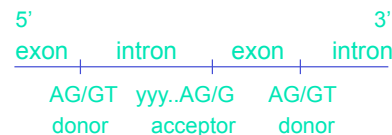
## Genes in Eukaryotes

As in prokaryotes (but maybe more variable)

- promoters
- start/stop transcription
- start/stop translation

New Features:

- polyA site/tail
- introns, exons, splicing
- branch point signal
- alternative splicing



42

## Characteristics of human genes

(Nature, 2/2001, Table 21)

	Median	Mean	Sample (size)
Internal exon	122 bp	145 bp	RefSeq alignments to draft genome sequence, with confirmed intron boundaries (43,317 exons)
Exon number	7	8.8	RefSeq alignments to finished seq (3,501 genes)
Introns	1,023 bp	3,365 bp	RefSeq alignments to finished seq (27,238 introns)
3' UTR	400 bp	770 bp	Confirmed by mRNA or EST on chromo 22 (689)
5' UTR	240 bp	300 bp	Confirmed by mRNA or EST on chromo 22 (463)
Coding seq	1,100 bp	1340bp	Selected RefSeq entries (1,804)*
(CDS)	367 aa	447 aa	
Genomic span	14 kb	27 kb	Selected RefSeq entries (1,804)*

\* 1,804 selected RefSeq entries were those with full-length unambiguous alignment to finished sequence

43



# Big Genes

Many genes are over 100 kb long,  
 Max known: dystrophin gene (DMD), 2.4 Mb.  
 The variation in the size distribution of coding sequences and exons is less extreme, although there are remarkable outliers.

The titin gene has the longest currently known coding sequence at 80,780 bp; it also has the largest number of exons (178) and longest single exon (17,106 bp).

RNApol rate: 1.2-2.5 kb/min =>16 hours to transcribe DMD

Nature 2/2001

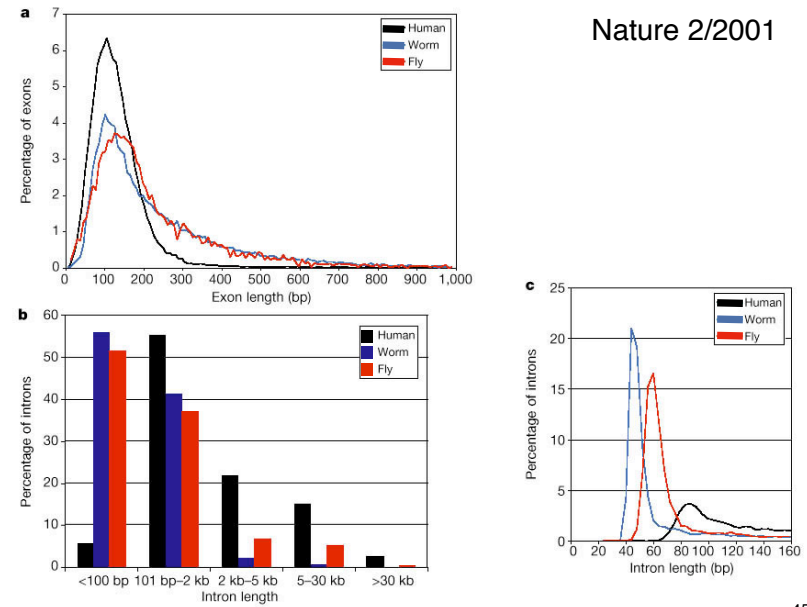
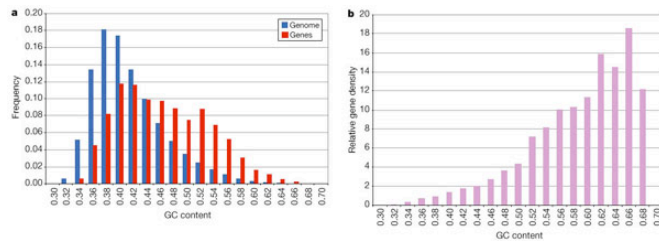


Figure 36 GC content

Nature 2/2001



**a:** Distribution of GC content in genes and in the genome. For 9,315 known genes mapped to the draft genome sequence, the local GC content was calculated in a window covering either the whole alignment or 20,000 bp centered on midpoint of the alignment, whichever was larger. Ns in the sequence were not counted. GC content for the genome was calculated for adjacent nonoverlapping 20,000-bp windows across the sequence. Both distributions normalized to sum to one.

**b:** Gene density as a function of GC content (= ratios of data in a. Less accurate at high GC because the denominator is small)

**c:** Dependence of mean exon and intron lengths on GC content. The local GC content, based on alignments to finished sequence only, calculated from windows covering the larger of feature size or 10,000 bp centered on it

# Computational Gene Finding?

How do we algorithmically account for all this complexity...

## A Case Study -- Genscan

C Burge, S Karlin (1997), "Prediction of complete gene structures in human genomic DNA", Journal of Molecular Biology , 268: 78-94.

48

## Training Data

238 multi-exon genes  
 142 single-exon genes  
 total of 1492 exons  
 total of 1254 introns  
 total of 2.5 Mb

NO alternate splicing, none > 30kb, ...

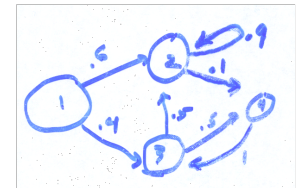
49

## Performance Comparison

Program	Accuracy						
	per nuc.		per exon				
	Sn	Sp	Sn	Sp	Avg.	ME	WE
GENSCAN	0.93	0.93	0.78	0.81	0.80	0.09	0.05
FGENEH	0.77	0.88	0.61	0.64	0.64	0.15	0.12
GeneID	0.63	0.81	0.44	0.46	0.45	0.28	0.24
Genie	0.76	0.77	0.55	0.48	0.51	0.17	0.33
GenLang	0.72	0.79	0.51	0.52	0.52	0.21	0.22
GeneParser2	0.66	0.79	0.35	0.40	0.37	0.34	0.17
GRAIL2	0.72	0.87	0.36	0.43	0.40	0.25	0.11
SORFIND	0.71	0.85	0.42	0.47	0.45	0.24	0.14
Xpound	0.61	0.87	0.15	0.18	0.17	0.33	0.13
GeneID#	0.91	0.91	0.73	0.70	0.71	0.07	0.13
GeneParser3	0.86	0.91	0.56	0.58	0.57	0.14	0.09

After Burge&Karlin, Table 1. Sensitivity, Sn = TP/AP; Specificity, Sp = TP/PP

## Generalized Hidden Markov Models



$\pi$ : Initial state distribution

$a_{ij}$ : Transition probabilities

One submodel per state

Outputs are *strings* gen'ed by submodel

Given length  $L$

Pick start state  $q_1$  ( $\sim \pi$ )

While  $\sum d_i < L$

Pick  $d_i$  & string  $s_i$  of length  $d_i$   $\sim$  submodel for  $q_i$

Pick next state  $q_{i+1}$  ( $\sim a_{ij}$ )

Output  $s_1 s_2 \dots$

51



## Decoding

A "parse"  $\phi$  of  $s = s_1s_2\dots s_L$  is a pair  $d = d_1d_2\dots d_k$ ,  $q = q_1q_2\dots q_k$  with  $\sum d_i = L$

A forward/backward-like alg calculates, e.g.:

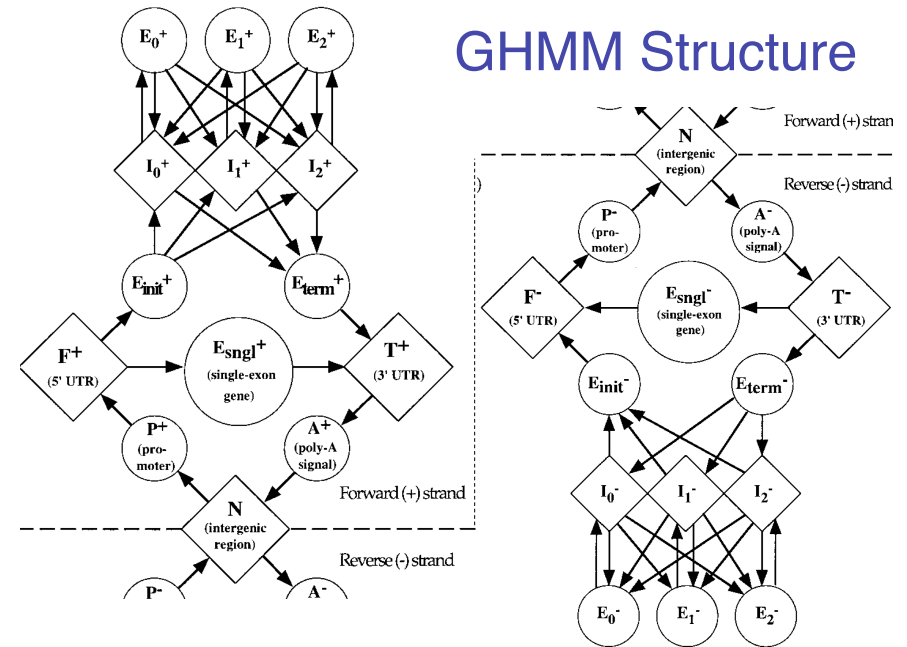
$$Pr(\text{generate } s_1s_2\dots s_i \text{ \& end in state } q_k)$$

(summing over possible predecessor states  $q_{k-1}$  and possible  $d_k$ , etc.)

$$Pr(\phi | s) = \frac{Pr(\phi, s)}{Pr(s)} \dots$$

52

## GHMM Structure



## Length Distributions

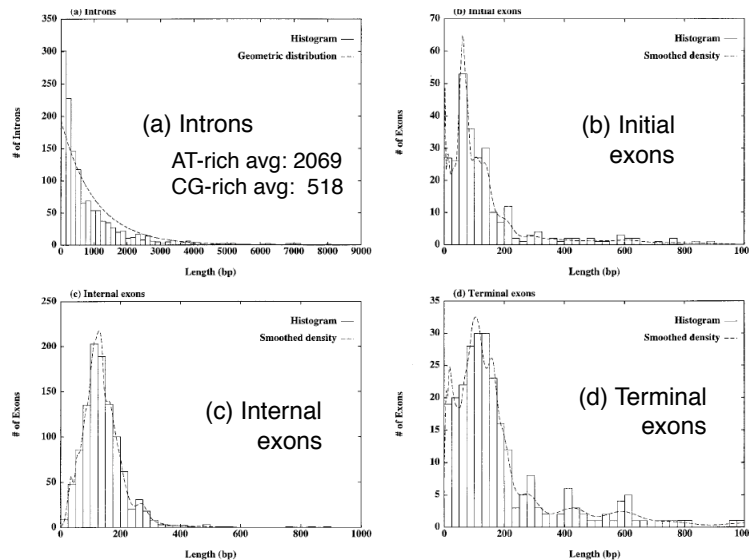


Figure 4. Length distributions are shown for (a) 1254 introns; (b) 238 initial exons; (c) 1151 internal exons; and (d) 238 terminal exons from the 238 multi-exon genes of the learning set  $\mathcal{L}$ . Histograms (continuous lines) were derived with a bin size of 300 bp in (a), and 25 bp in (b), (c), (d). The broken line in (a) shows a geometric (exponential) distribution with parameters derived from the mean of the intron lengths; broken lines in (b), (c) and (d) are the smoothed empirical distributions of exon lengths used by GENSCAN (details given by Burge, 1997). Note different horizontal and vertical scales are used in (a), (b), (c), (d) and that multimodality in (b) and (d) may, in part, reflect relatively small sample sizes.

## Effect of G+C Content

Group	I	II	III	IV
C ≠ G% range	<43	43-51	51-57	>57
Number of genes	65	115	99	101
Est. proportion single-exon genes	0.16	0.19	0.23	0.16
Codelen: single-exon genes (bp)	1130	1251	1304	1137
Codelen: multi-exon genes (bp)	902	908	1118	1165
Introns per multi-exon gene	5.1	4.9	5.5	5.6
Mean intron length (bp)	2069	1086	801	518
Est. mean transcript length (bp)	10866	6504	5781	4833
Isochore	L1+L2	H1+H2	H3	H3
DNA amount in genome (Mb)	2074	1054	102	68
Estimated gene number	22100	24700	9100	9100
Est. mean intergenic length	83000	36000	5400	2600
<b>Initial probabilities:</b>				
Intergenic (N)	0.892	0.867	0.54	0.418
Intron (I+, I-)	0.095	0.103	0.338	0.388
5' Untranslated region (F+, F-)	0.008	0.018	0.077	0.122
3' Untranslated region (T+, T-)	0.005	0.011	0.045	0.072

55

## Submodels

### 5' UTR

L ~ geometric(769 bp), s ~ MM(5)

### 3' UTR

L ~ geometric(457 bp), s ~ MM(5)

### Intergenic

L ~ geometric(GC-dependent), s ~ MM(5)

### Introns

L ~ geometric(GC-dependent), s ~ MM(5)

56

## Submodel: Exons

Inhomogeneous 3-periodic 5th order Markov models

Separate models for low GC (<43%), high GC

Track “phase” of exons, i.e. reading frame.

57

## Signal Models I: WMM's

### Polyadenylation

6 bp, consensus AATAAA

### Translation Start

12 bp, starting 6 bp before start codon

### Translation stop

A stop codon, then 3 bp WMM

58

## Signal Models II: more WMM's

### Promoter

70% TATA

15 bp TATA WMM

s ~ null, L ~ Unif(14-20)

8 bp cap signal WMM

30% TATA-less

40 bp null

59

## Signal Models III: W/WAM's

### Acceptor Splice Site (3' end of intron)

[-20..+3] relative to splice site modeled by "1st order weight array model"

### Branch point & polypyrimidine tract

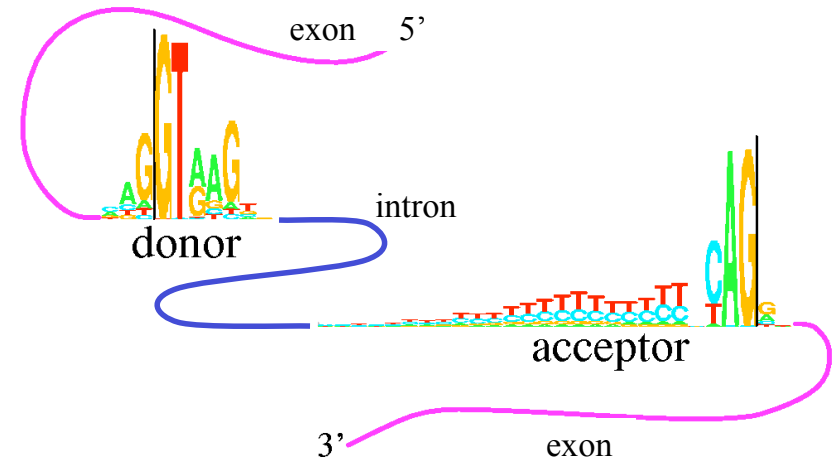
Hard. Even weak consensus like YYRAY found in [-40..-21] in only 30% of training

"Windowed WAM": 2nd order WAM, but averaged over 5 preceding positions

"captures weak but detectable tendency toward YYY triplets and certain branch point related triplets like TGA, TAA, ..."

60

## What's in the Primary Sequence?



## Signal Models IV: Maximum Dependence Decomposition

Donor splice sites (5' end of intron) show dependencies between non-adjacent positions, e.g. poor match at one end compensated by strong match at other end, 6 bp away

Model is basically a decision tree

Uses  $\chi^2$  test to quantitate dependence

62

## Are A & B independent ?

	B	not B	
A	8	4	12
not A	2	6	8
	10	10	20

$$\chi^2 = \sum_i \frac{(\text{observed}_i - \text{expected}_i)^2}{\text{expected}_i}$$

"Expected" means expected assuming independence, e.g. expect B 10/20; A 12/20; both 120/400\*20 = 6, etc.

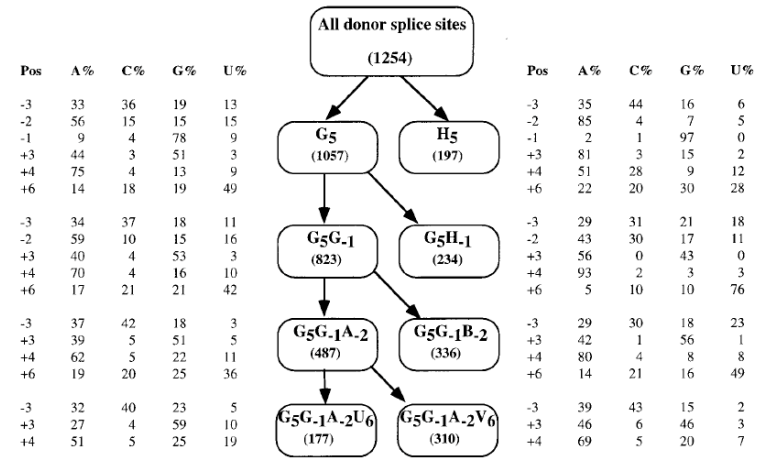
Look up in table (or approximate as normal).

63

# $\chi^2$ test for independence

i	Con	j:	-3	-2	-1	+3	+4	+5	+6	Sum
-3	c/a	---	61.8*	14.9	5.8	20.2*	11.2	18.0*	131.8*	
-2	A	115.6*	---	40.5*	20.3*	57.5*	59.7*	42.9*	336.5*	
-1	G	15.4	82.8*	---	13.0	61.5*	41.4*	96.6*	310.8*	
+3	a/g	8.6	17.5*	13.1	---	19.3*	1.8	0.1	60.5*	
+4	A	21.8*	56.0*	62.1*	64.1*	---	56.8*	0.2	260.9*	
+5	G	11.6	60.1*	41.9*	93.6*	146.6*	---	33.6*	387.3*	
+6	t	22.2*	40.7*	103.8*	26.5*	17.8*	32.6*	---	243.6*	

\* means chi-squared p-value < .001

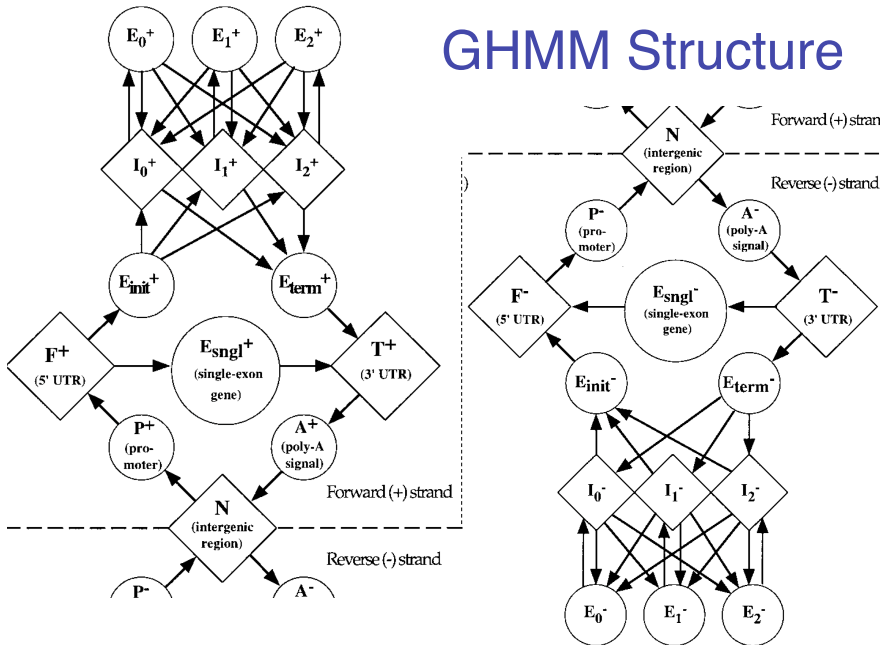


All sites: Position

Base	-3	-2	-1	+1	+2	+3	+4	+5	+6
A%	33	60	8	0	0	49	71	6	15
C%	37	13	4	0	3	7	5	19	
G%	18	14	81	100	0	45	12	84	20
U%	12	13	7	0	100	3	9	5	46

U1 snRNA: 3' G U C C A U U C A 5'

# GHMM Structure



# Summary of Burge & Karlin

Coding DNA & control signals nonrandom

Weight matrices, WAMs, etc. for controls

Codon frequency, etc. for coding

GHMM nice for overall architecture

Careful attention to small details pays

## Problems with BK training set

- 1 gene per sequence
- Annotation errors
- Single exon genes over-represented?
- Highly expressed genes over-represented?
- Moderate sized genes over-represented?  
(none > 30 kb) ...
- Similar problems with other training sets, too

68

## Problems with all methods

- Pseudo genes
- Short ORFs
- Sequencing errors
- Non-coding RNA genes & spliced UTR's
- Overlapping genes
- Alternative splicing/polyadenylation
- Hard to find novel stuff – not in training
- Species-specific weirdness – spliced leaders, polycistronic transcripts, RNA editing...

69

## Other important ideas

Database search - does gene you're predicting look anything like a known protein?

Comparative genomics - what does this region look like in related organisms?

70