

# Genome Characteristics and Annotation

---

COMP 571 - Spring 2015  
Luay Nakhleh, Rice University

# Outline

- \* Gene prediction in prokaryotic genomes
- \* Features used in eukaryotic gene detection
- \* Predicting eukaryotic gene signals
- \* Complete eukaryotic gene models
- \* Genome annotation

# Gene Prediction in Prokaryotic Genomes

# A simple gene structure

- \* Although introns do exist in prokaryotes, they are extremely rare and often ignored by gene prediction tools.
- \* The relative simplicity of bacterial gene structure has led to some very successful gene prediction techniques that use functional signals, such as the ribosome-binding site, the stop codon that signals the end of translation, and other well-defined features.

# Illustration



- \* One can easily enumerate all potential open reading frames (ORFs) present in the genome.
- \* The longer the potential ORF, the more likely it is to really be a gene.
- \* A key problem then is to distinguish the true and false genes in the set of short potential ORFs of, say, 150 bases or fewer.

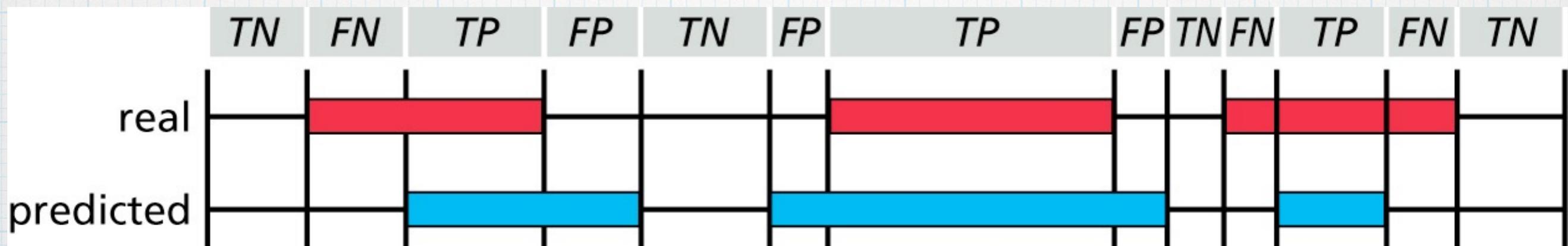
- \* A consequence of this situation is that by accepting some false positives, a gene detection method can achieve very high rates of detection!
- \* Put another way, these methods should really be detecting the false ORFs.

\* One must be wary of some of the high success rates quoted (even over 98%), and false positive rates would be more informative, but are often not quoted.

# Measures of Gene Prediction Accuracy

- \* In the field of gene prediction, accuracy can be measured at three different levels:
  - \* Coding nucleotides: The base level
  - \* Exon structure: The exon level
  - \* Protein product: The protein level

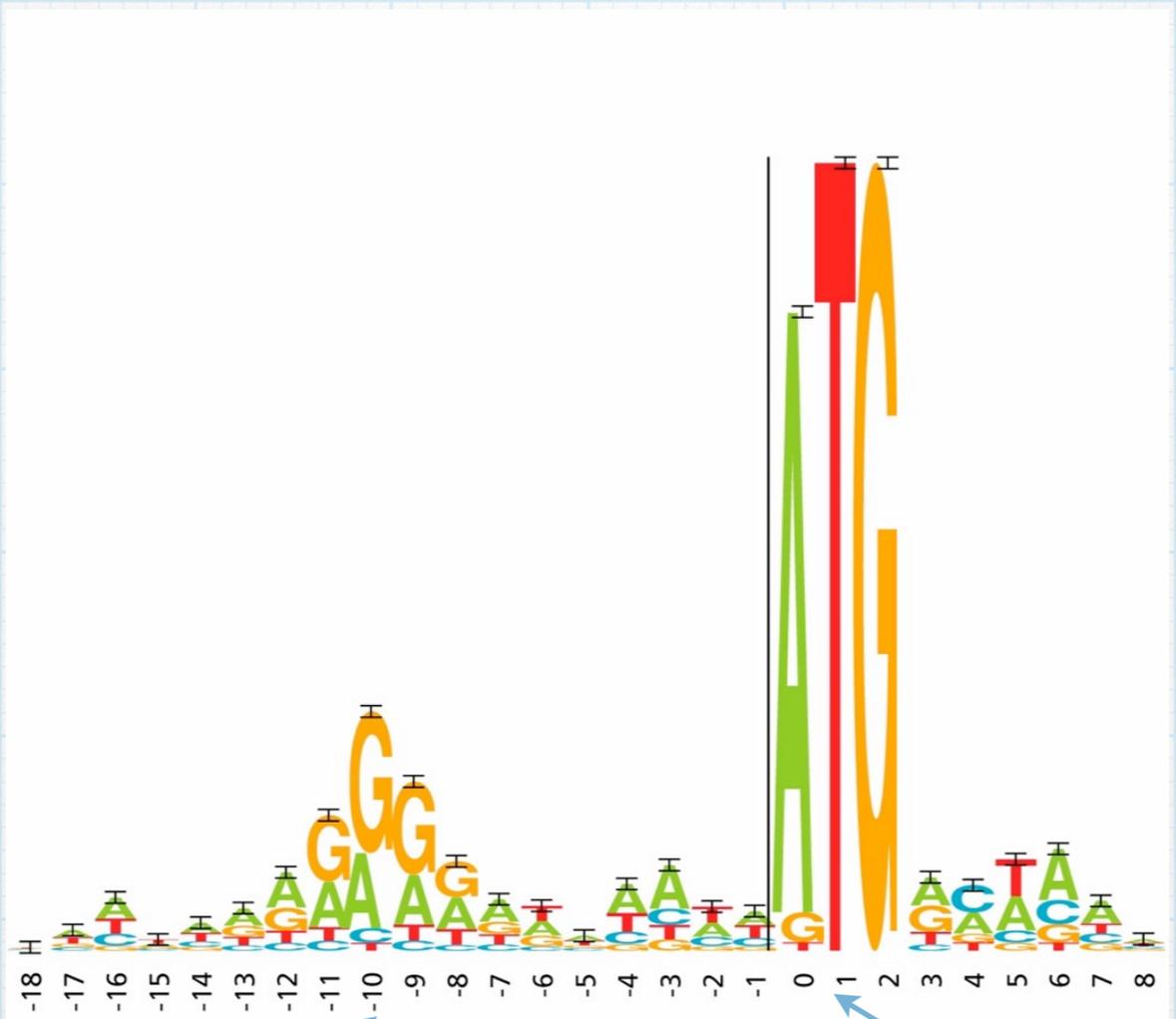
# Measures of Gene Prediction Accuracy at the Nucleotide Level



$$S_n = \frac{TP}{TP + FN}$$

$$S_p = \frac{TP}{TP + FP}$$

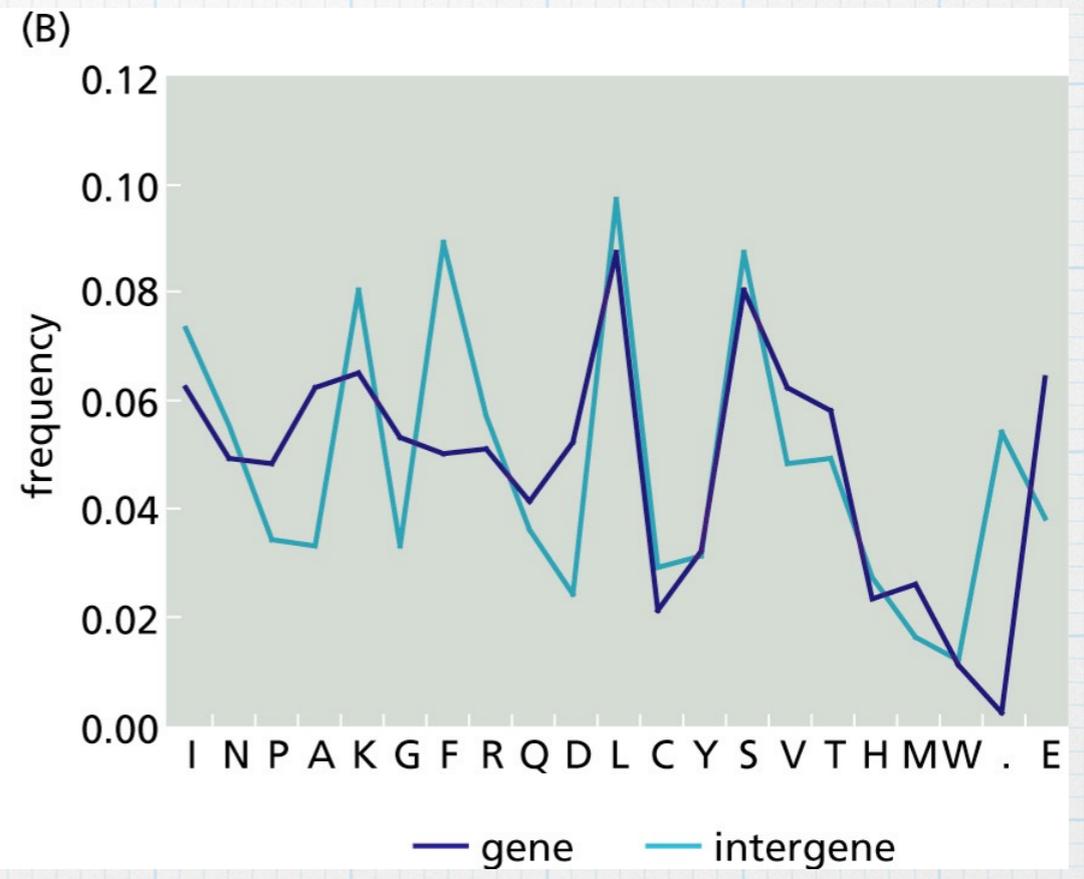
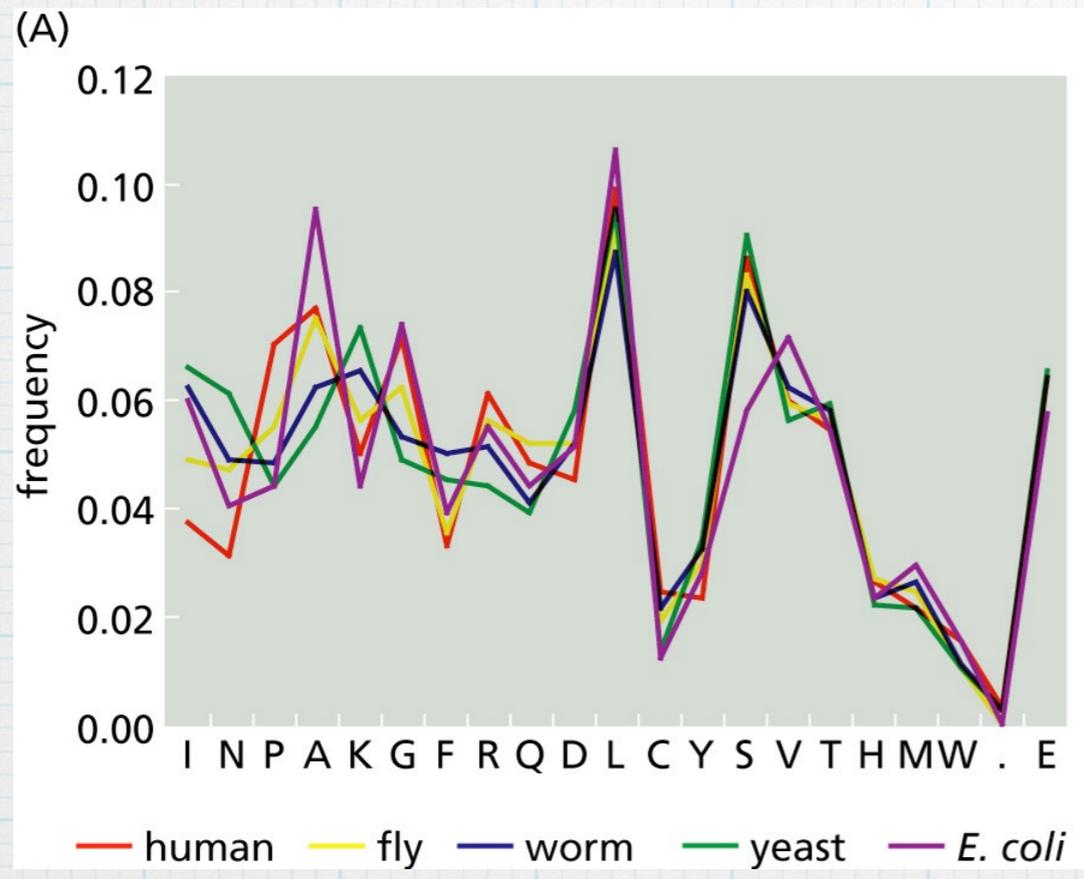
- \* The most basic characteristic of a gene is that it must contain an open reading frame (ORF) that begins with a start codon (ATG) and ends with a stop codon (TAA, TAG, or TGA).
- \* There are some exceptions (for example, E. coli uses GTG for 9% and TTG for 0.5% of start codons).



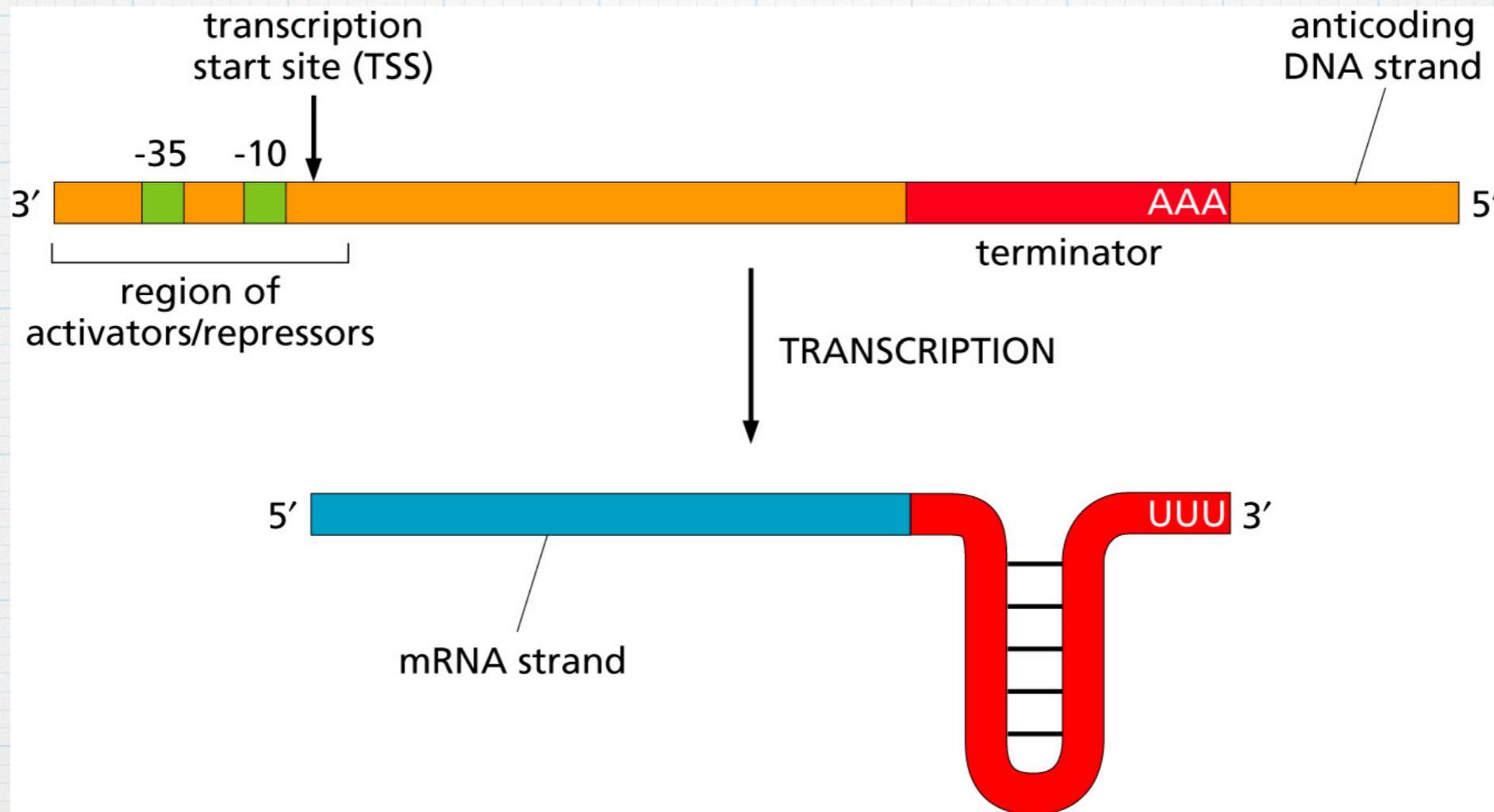
Ribosome-binding site

Start codon in E. coli genes

**\* Another characteristic that can be used to detect genes is the relative frequency of codon occurrences.**



# Gene Structure in Prokaryotes



\* Bacterial promoters typically occur immediately before the position of the **transcription start site (TSS)**, and contain two characteristic short sequences, or motifs, that are almost the same in the promoters for different genes.

\* The termination of transcription is controlled by the **terminator signal** which in bacteria differs from the promoter is that it is active when transcribed to form the end of the mRNA strand (forms a loop structure that prevents the transcription apparatus from continuing).

\* Single type of RNA polymerase transcribes all genes.

# Algorithms for Gene Detection in Prokaryotes

- \* GeneMark
- \* GeneMark.hmm
- \* GLIMMER
- \* ORPHEUS
- \* ...

# GeneMark

- \* GeneMark uses a fifth-order Markov chain model to represent the statistics of coding and noncoding reading frames.
- \* The method uses the dicodon statistics to identify coding regions.

# GeneMark

$$P(a|b_1b_2b_3b_4b_5) = \frac{n_{b_1b_2b_3b_4b_5a}}{\sum_{\alpha \in \{A,C,T,G\}} n_{b_1b_2b_3b_4b_5\alpha}}$$


The number of times  $b_1b_2b_3b_4b_5\alpha$  occurs in the training data

GeneMark assumes each reading frame has unique dicodon statistics, and thus has its own model probabilities  $P_1, P_2, P_3, P_4, P_5, P_6$ .

For noncoding regions, there is  $P_{nc}(a|b_1b_2b_3b_4b_5)$ .

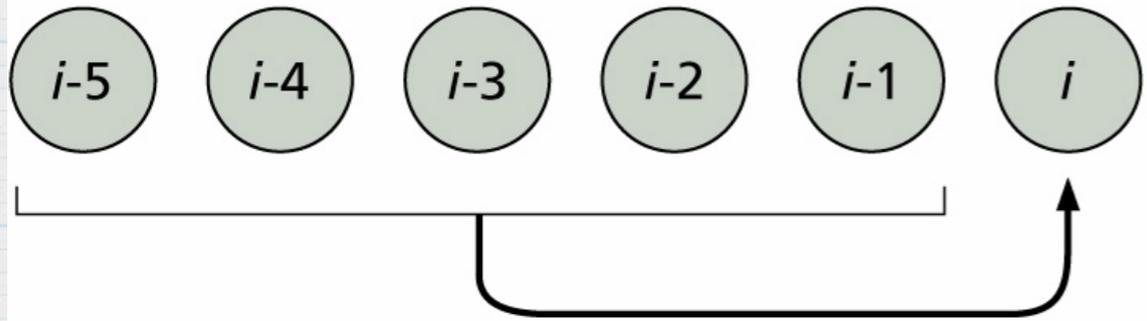
# GeneMark

- \* For example, the probability of obtaining a sequence  $x = x_1 x_2 \dots x_9$  if  $x_1 x_2 x_3$  is a translated codon (that is,  $x_9$  is in the third position of a translated codon) is given by

$$\mathbf{P}(x|3) = \mathbf{P}_2(x_1 x_2 x_3 x_4 x_5) \mathbf{P}_2(x_6 | x_1 x_2 x_3 x_4 x_5) \mathbf{P}_3(x_7 | x_2 x_3 x_4 x_5 x_6) \\ \times \mathbf{P}_1(x_8 | x_3 x_4 x_5 x_6 x_7) \mathbf{P}_2(x_9 | x_4 x_5 x_6 x_7 x_8)$$

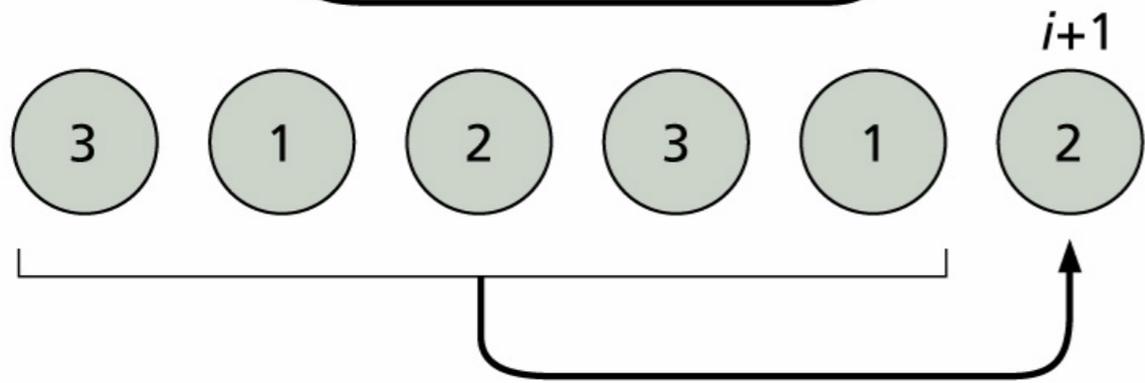
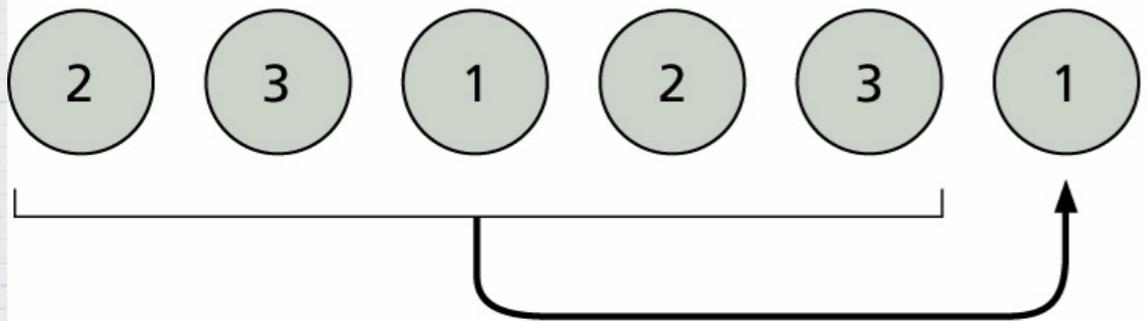
**This is called a periodic, phased, or inhomogeneous Markov model.**

(A)

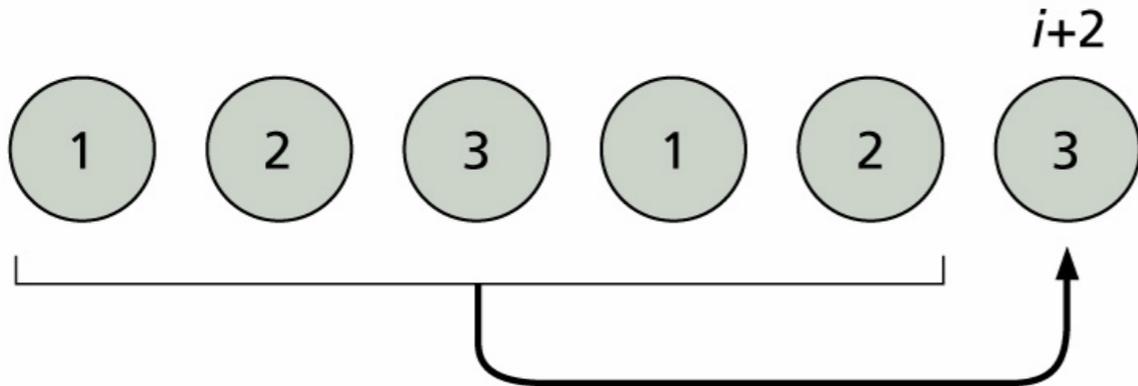


homogeneous

(B)



inhomogeneous



# GeneMark

- \* We want  $P(3|x)$ , which can be derived using Bayes rule as

$$P(3|x) = \frac{P(x|3)P(3)}{P(x|nc)P(nc) + \sum_{i=1}^6 P(x|i)P(i)}$$

- \* Similar formulas can be derived for  $P(i|x)$  for any value of  $i$ .

# GeneMark

- \* In GeneMark,  $P(nc)$  was assumed to be  $1/2$ , and  $P(1)$ - $P(6)$  were assumed to all be  $1/12$ .
- \* Sliding windows of 96 nucleotides were scored in steps of 12 nucleotides.
- \* If  $P(i|x)$  exceeds a certain threshold, the window is predicted to be in coding reading frame  $i$ .

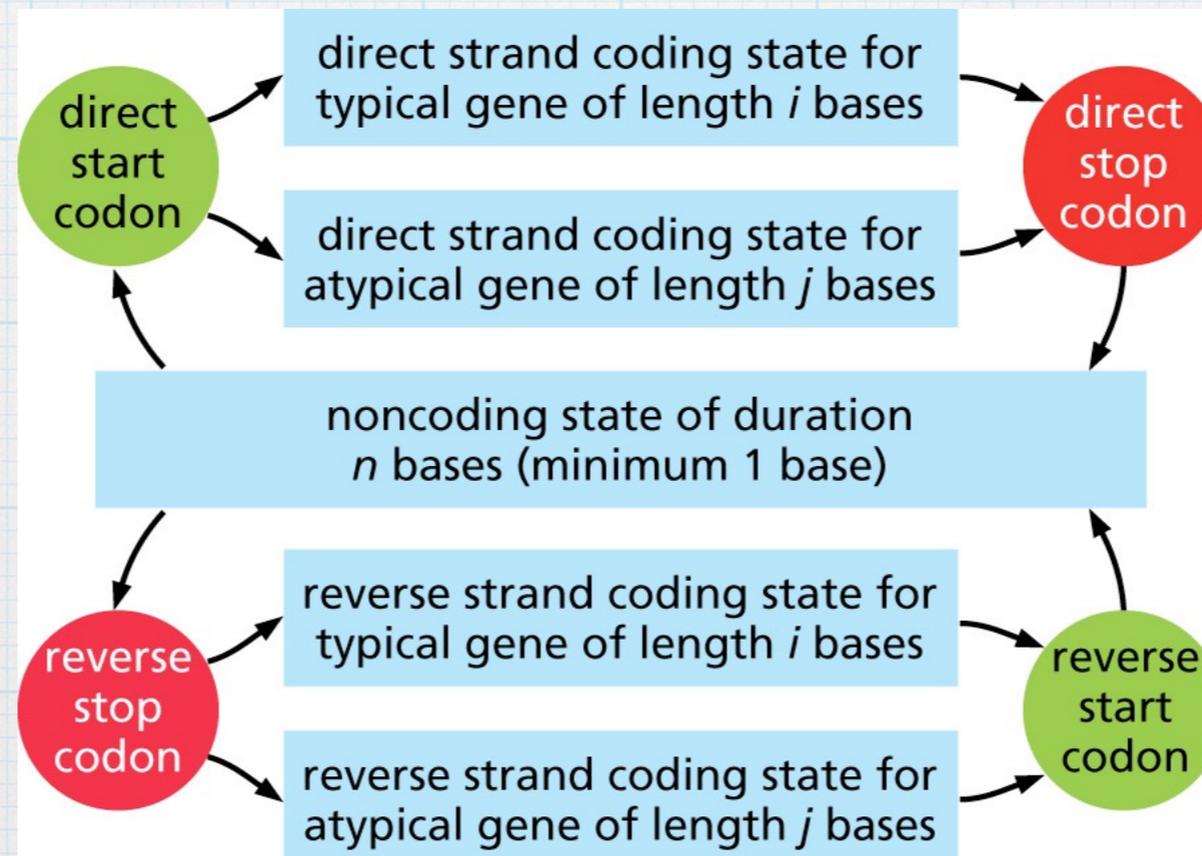
# GeneMark

- \* The final predicted gene boundaries are defined by start and stop codons in that reading frame.

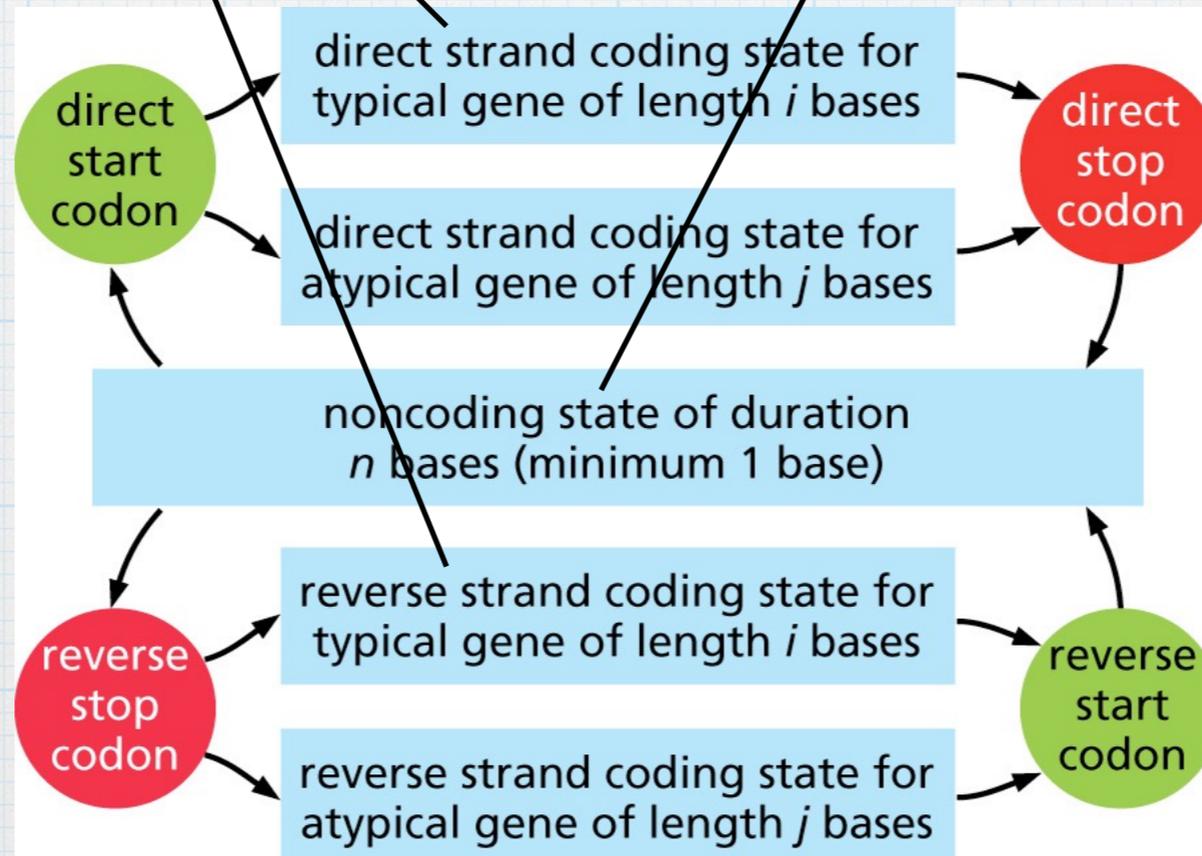
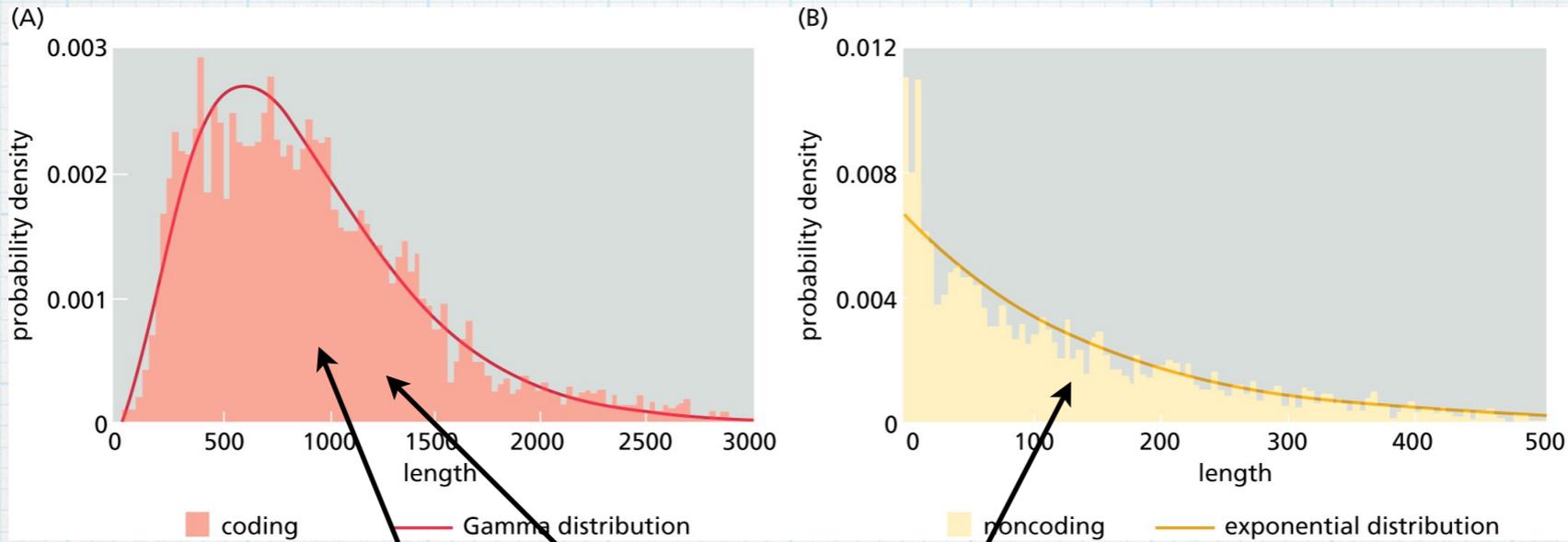
# GeneMark.hmm

- \* GeneMark uses a sliding window, and doesn't do a good job at defining the gene boundaries.
- \* GeneMark.hmm is an extension to ameliorate these issues.

# GeneMark.hmm



# GeneMark.hmm



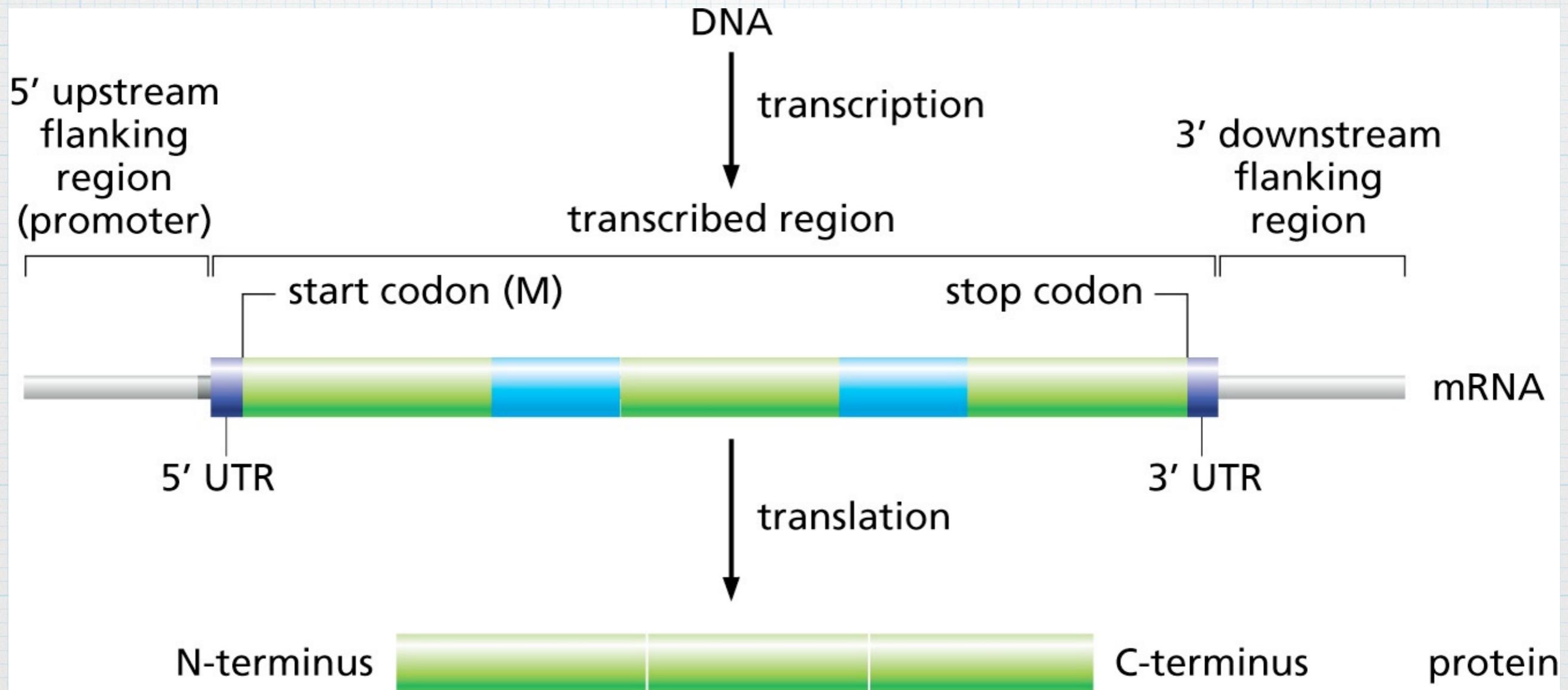
# Features Used in Eukaryotic Gene Detection

- \* Many of the principles that apply to the detection of genes in prokaryotes also apply to gene finding in eukaryotes.
- \* For example, the coding regions of eukaryotic genomes have distinct base statistics similar to those found in prokaryotes.

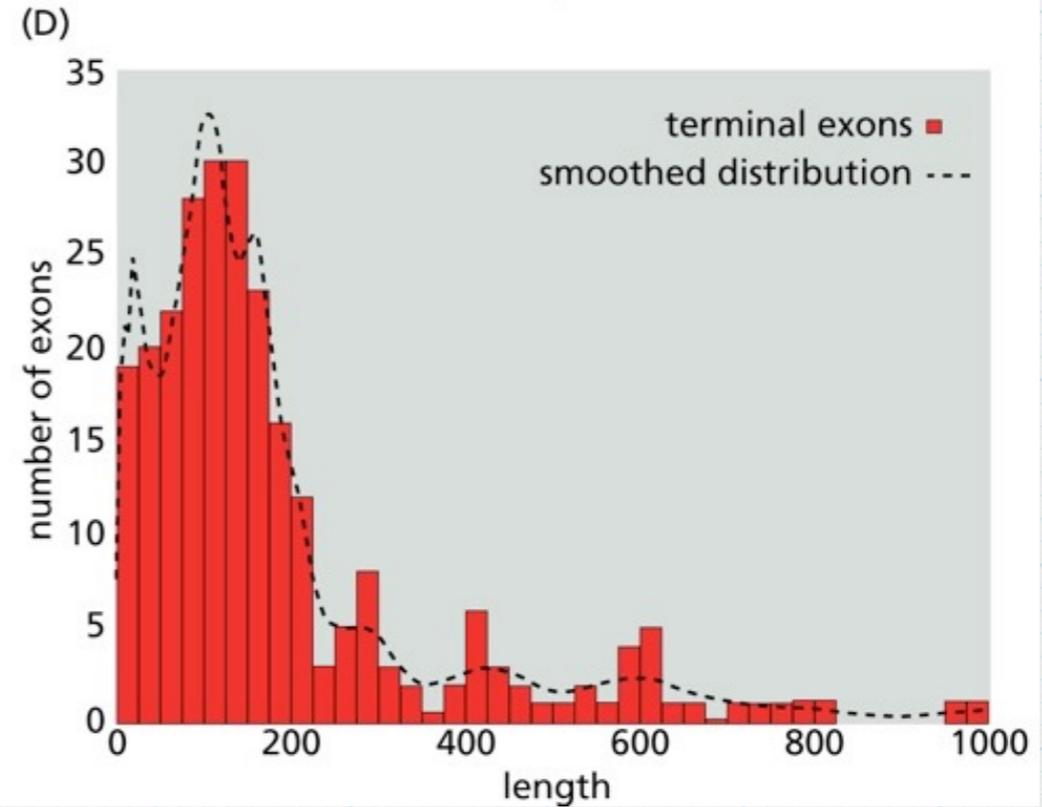
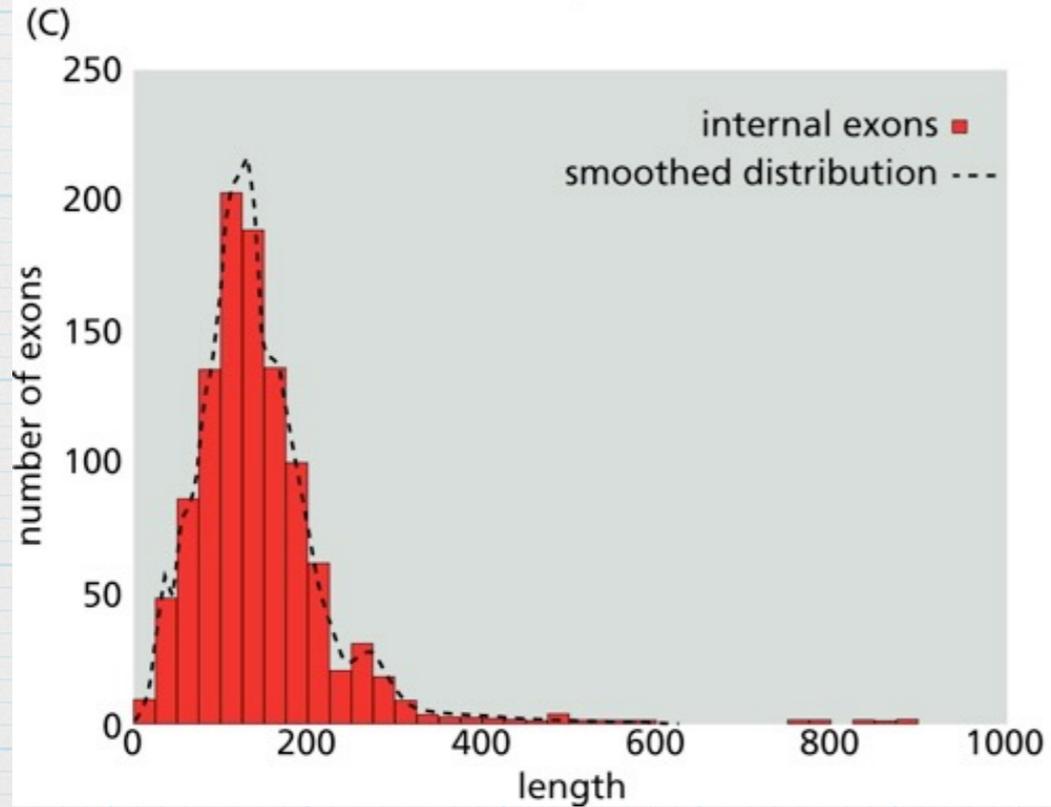
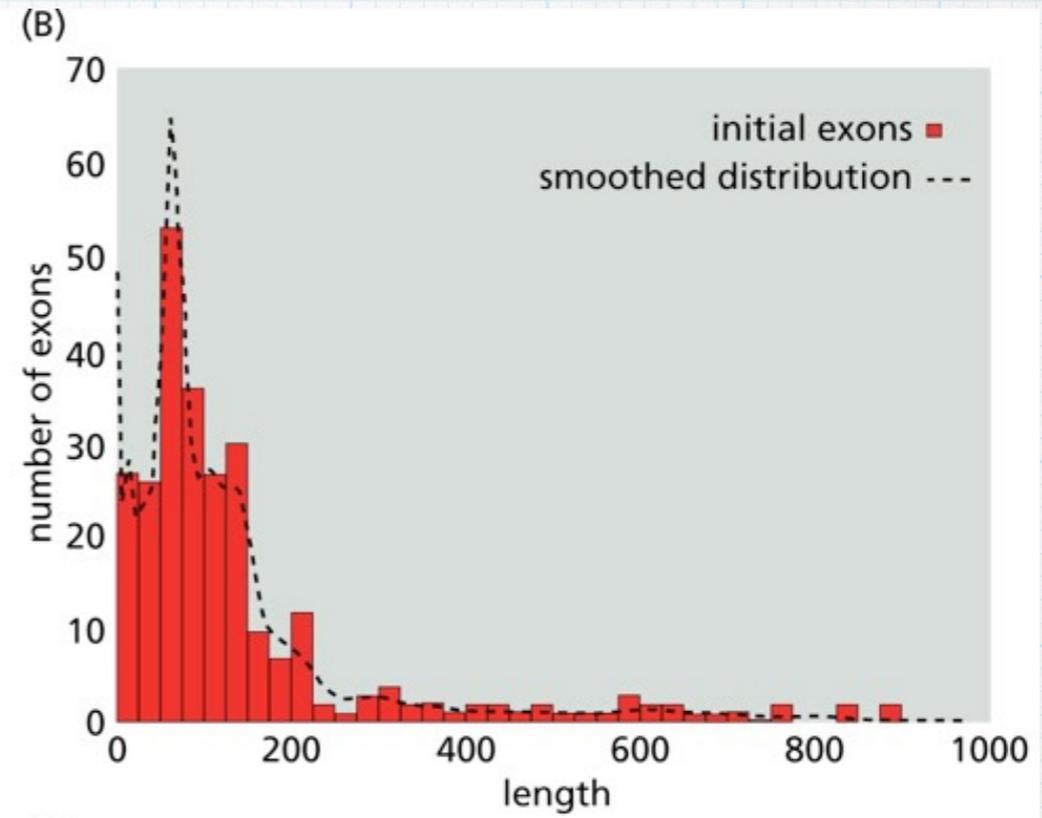
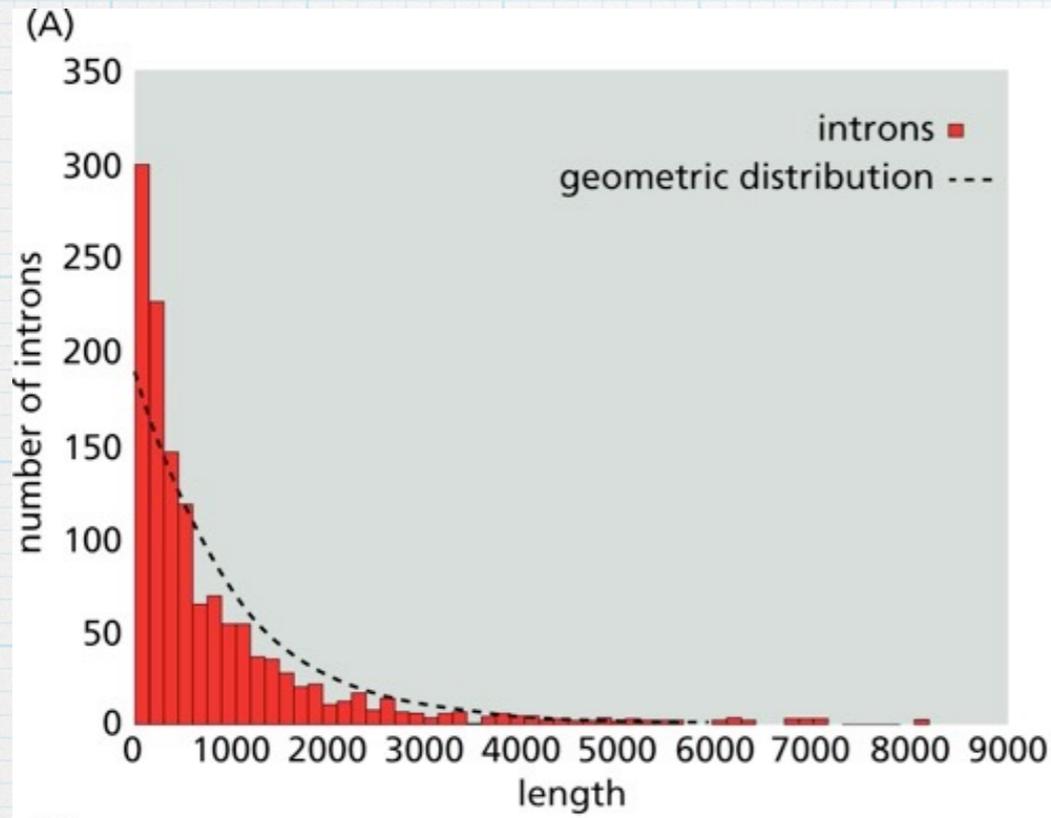
- \* In addition, although the signals differ, there are equivalent transcription and translation start and stop signals.**

- \* A crucial difference in gene structure causes eukaryotic gene detection to be far harder: there are numerous introns present in many genes.

# From Eukaryotic DNA to Protein



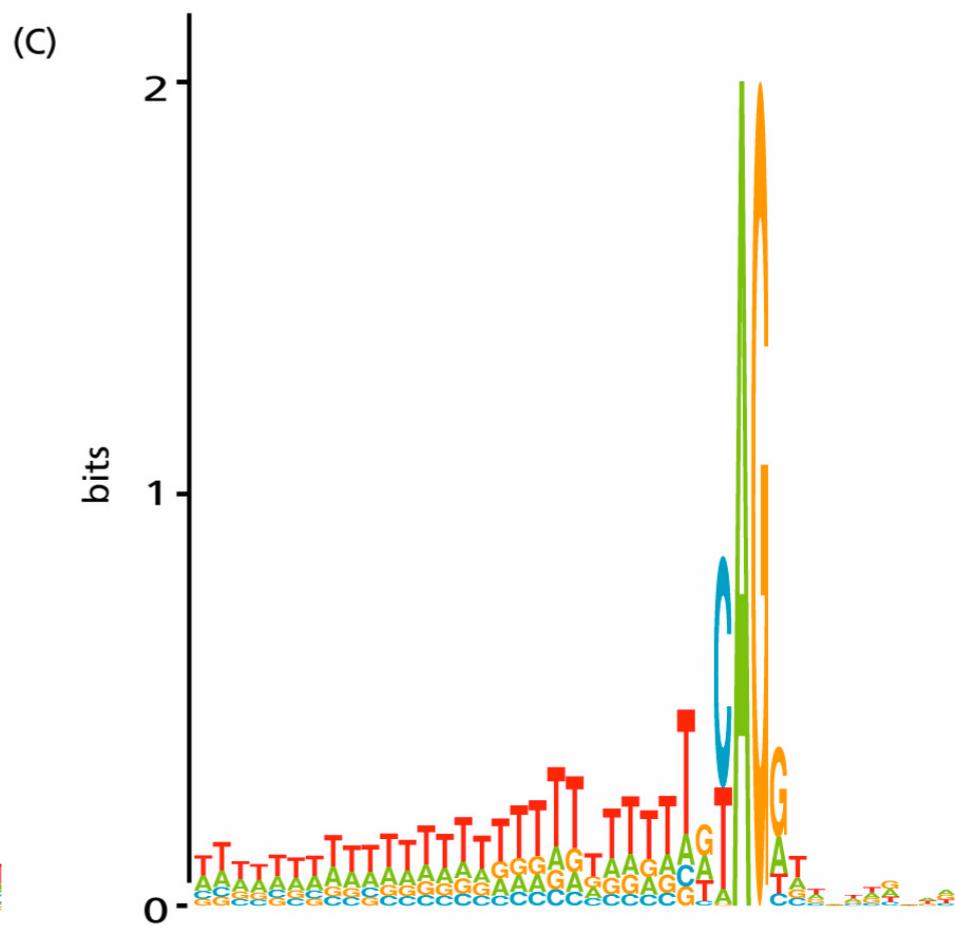
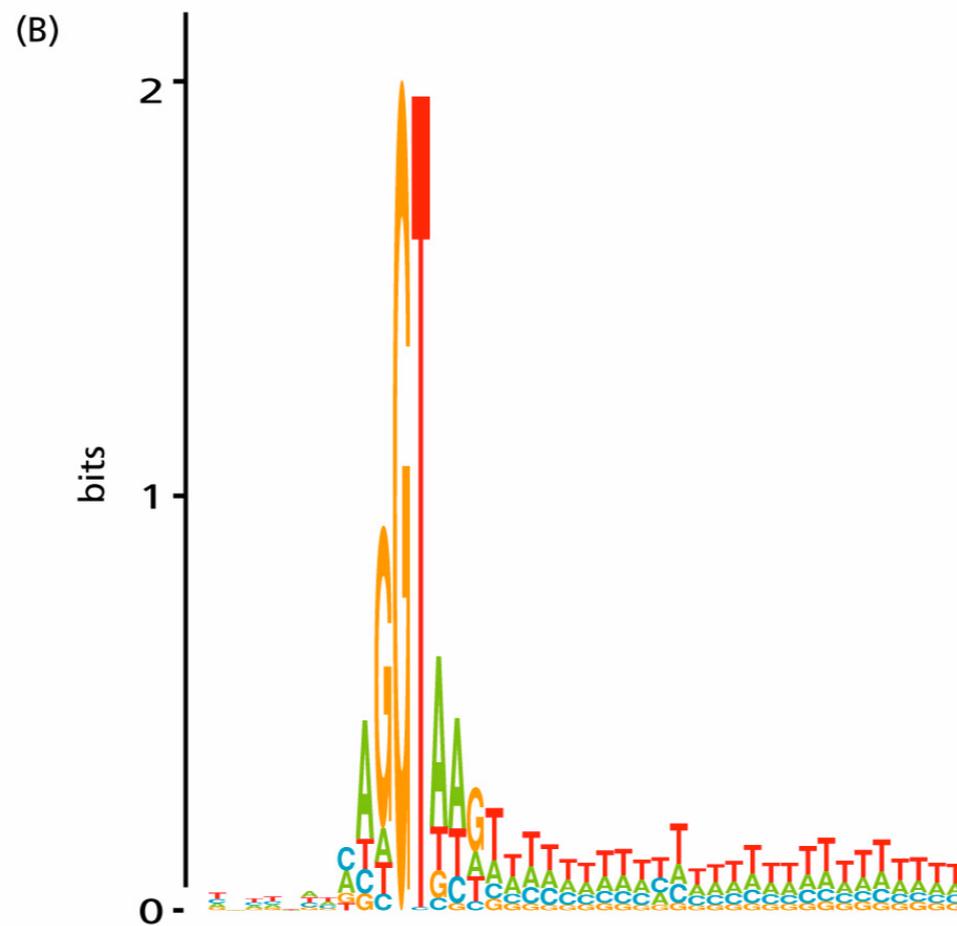
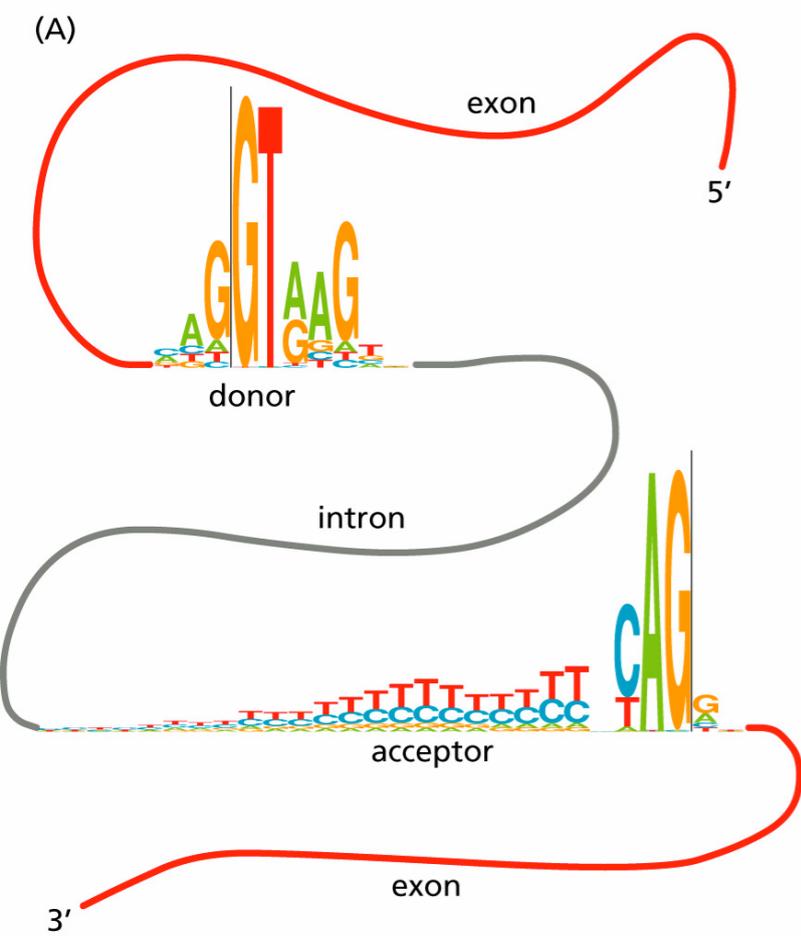
- \* The length of the protein-coding segments (exons) is on average smaller in eukaryotes than in prokaryotes, resulting in poorer base statistics, and making their detection more difficult.



# Distributions in the human genome

- \* An additional difference that can also cause difficulties is that the density of genes in most segments of eukaryotic genomes is significantly less than in prokaryotes.

- \* The splice signals at intron-exon boundaries are quite variable, making them hard to locate accurately.



human donor/acceptor sites

donor sites

acceptor sites

in Arabidopsis

# Alternative Splicing

- \* A particularly difficult problem can arise in eukaryotic genomes when moving from gene detection to protein prediction, a trivial step in prokaryotes.
- \* The splicing of introns in the RNA is not always identical for a given gene (the phenomenon of alternative splicing).

# Alternative Splicing

- \* Alternative splicing can give rise to the production of two or more different proteins from the same gene, and these are often known as splice variants.

# Promoter Sequences and Binding Sites for Transcription Factors

- \* A further difference between prokaryotic and eukaryotic gene structures is that the sequence signals in the upstream regions are much more variable in eukaryotes, both in composition and position.

# Promoter Sequences and Binding Sites for Transcription Factors

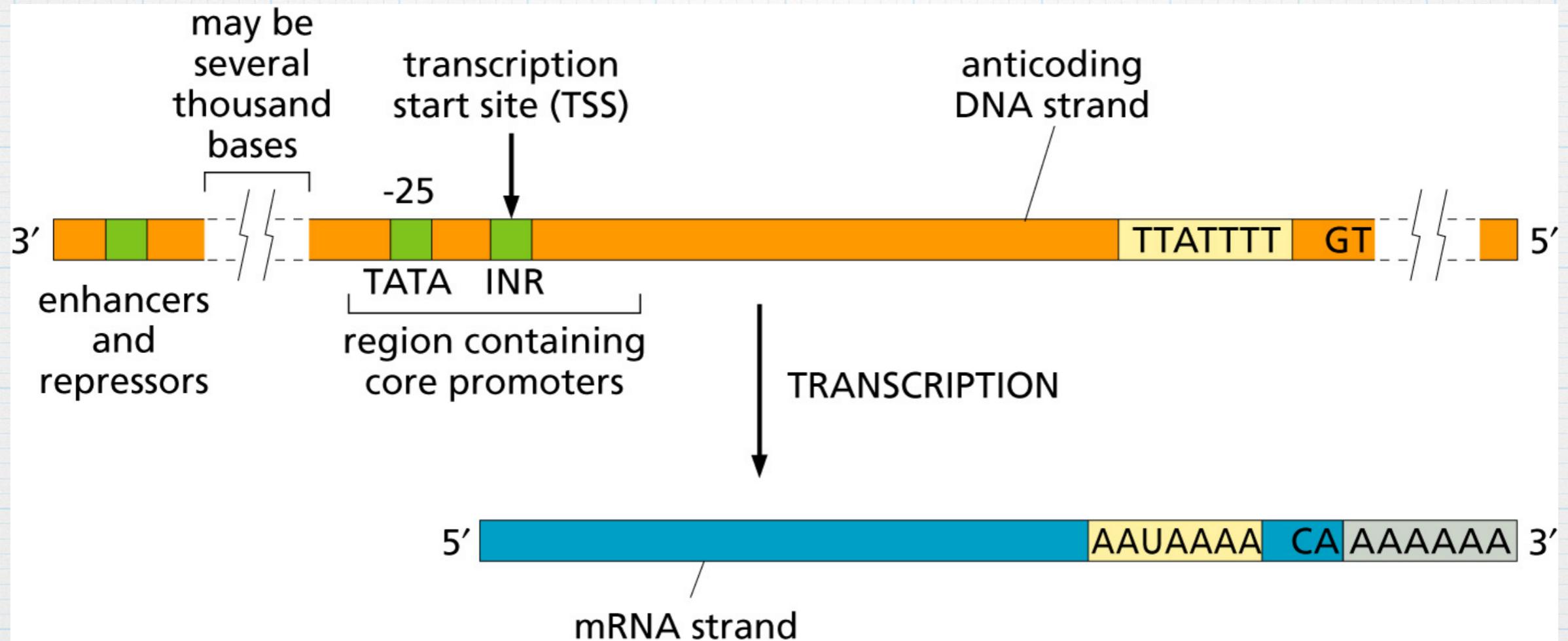
- \* The control of gene expression is more complex in eukaryotes than prokaryotes, and can be affected by many molecules binding the DNA in the region of the gene.

# Promoter Sequences and Binding Sites for Transcription Factors

- \* This leads to many more potential promoter binding signals spread over a much larger region (possibly several thousand bases) in the vicinity of the transcription start site.

# Predicting Eukaryotic Gene Signals

# Gene Structure in Eukaryotes



- \* Regulatory elements in eukaryotes are more complex.
- \* Three types of RNA polymerase transcribe genes: RNA polymerase II transcribes all protein coding genes, whereas other RNA polymerase types transcribe genes for tRNAs, rRNAs and other types of RNA

- \* In 1990 P. Bucher derived weight matrices to identify four separate RNA polymerase II promoter elements: the TATA box, the cap signal (INR), the CCATT box, and the GC box.

- \* Using more than 500 aligned eukaryotic sequences, the weights of different bases  $a$  at position  $u$  in a signal sequence were obtained from the general equation

$$w_u(a) = \ln \left( \frac{n_u(a)}{e_u(a)} + \frac{c}{100} \right) + c_u$$

number of occurrences of base  $a$  at position  $u$

expected number of bases  $a$  at position  $u$

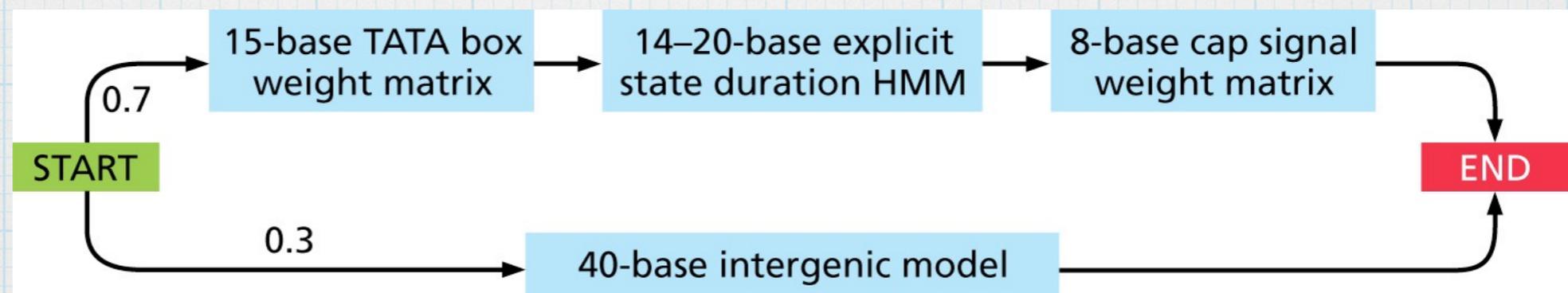
a small number (often 2)

adjusted to make the greatest  $w_u(a)$  zero



\* In GenScan (a popular gene detection method; more later), the promoter detection component uses Bucher's TATA-box and cap-signal models.

\* To avoid missing genes that lack a TATA-box, the model allows for both possibilities.



# Predicting Exons and Introns

- \* All internal introns and exons in a eukaryotic gene are delimited by the splice sites at which introns are cut out of the RNA transcript and the exon sequences joined together.

# Predicting Exons and Introns

- \* The splice sites have distinct sequence signals.
- \* There are programs that predict introns and exons without reference to splice sites, and other programs that predict splice sites without information about introns and exons.

# Predicting Exons and Introns

- \* For example, GenScan identifies eukaryotic coding regions by dicodon statistics, as in the prokaryotic example given earlier, but it uses an explicit state duration HMM based on the observed length distribution of real exons.
- \* The length of the potential exon is generated from this distribution, and its sequence generated with probabilities based on the dicodon statistics.

# Predicting Exons and Introns

- \* Measures of gene prediction accuracy at the exon level make use of:
  - \* AE: the number of actual (real) exons in the data
  - \* PE: The number of predicted exons
  - \* CE: the number of exons predicted exactly
  - \* ME: the number of real exons that do not overlap any of the predicted ones
  - \* WE: the number of predicted exons that have no overlap with the actual ones

# Predicting Exons and Introns

\* We then define:

\*  $S_{n1} = CE/AE$

$$S_{n2} = ME/AE$$

\*  $S_{p1} = CE/PE$

$$S_{p2} = WE/PE$$

# Predicting Exons and Introns

- \* Many of these ab initio prediction programs have been modified to take account of homology to experimental gene sequences.
- \* For example, GenScan allows BLAST hits to the genome sequence to be used as an extra guide.

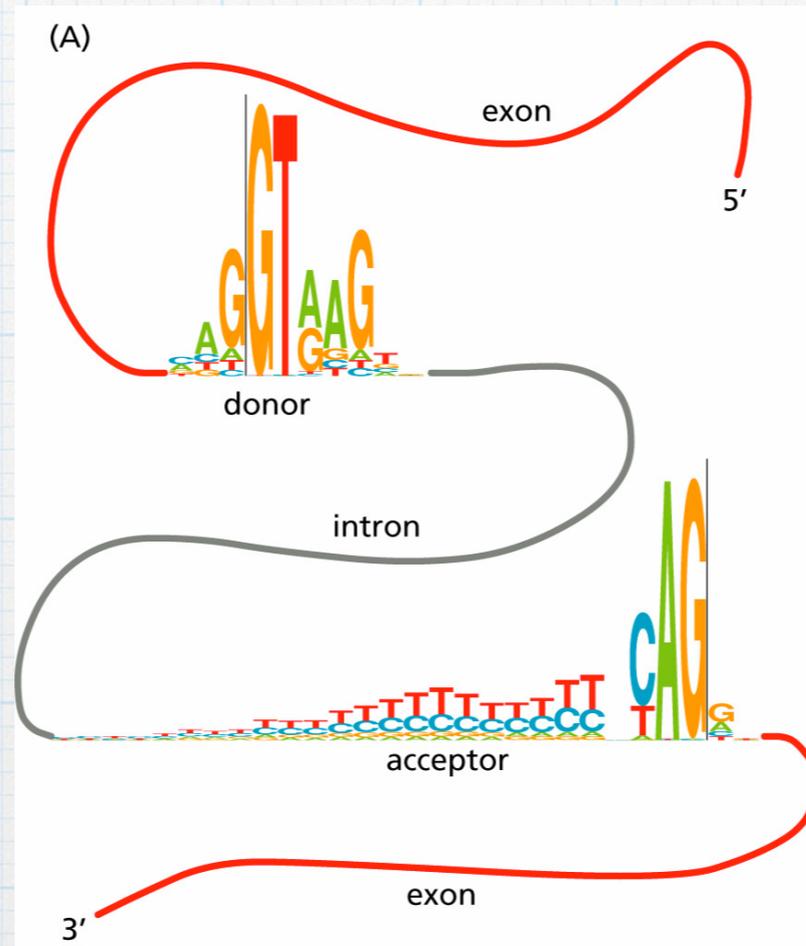
# Splice Site Prediction

- \* Most introns start with (i.e., have a 5' end) a GU dinucleotide in the RNA (GT in the DNA) at what is referred to as the donor splice site.
- \* The 3' end of introns (acceptor splice site) is mostly AG dinucleotide.

# Splice Site Prediction

- \* Thus, locating occurrences of AG and GT would identify all possible splice sites with these sequence properties, but in addition there would be about 30 to 100 false predicted sites for every true one.
- \* Properties of the surrounding sequence is used to reduce the false-negative prediction rate to a manageable level.

# Splice Site Prediction



# Splice Site Prediction

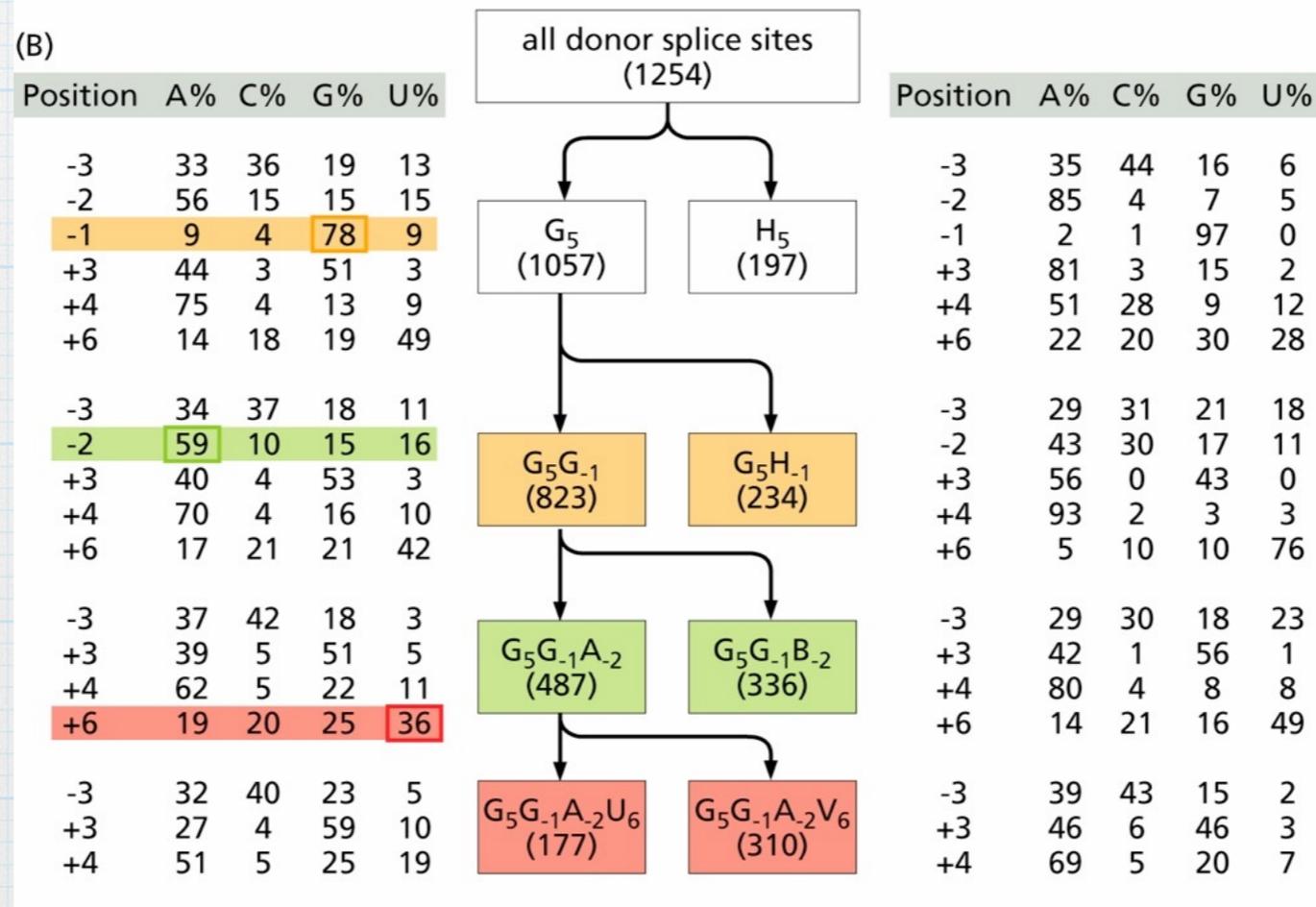
- \* GenScan uses a sophisticated variant of weighted matrices to predict donor splice sites.

# Splice Site Prediction

(A)

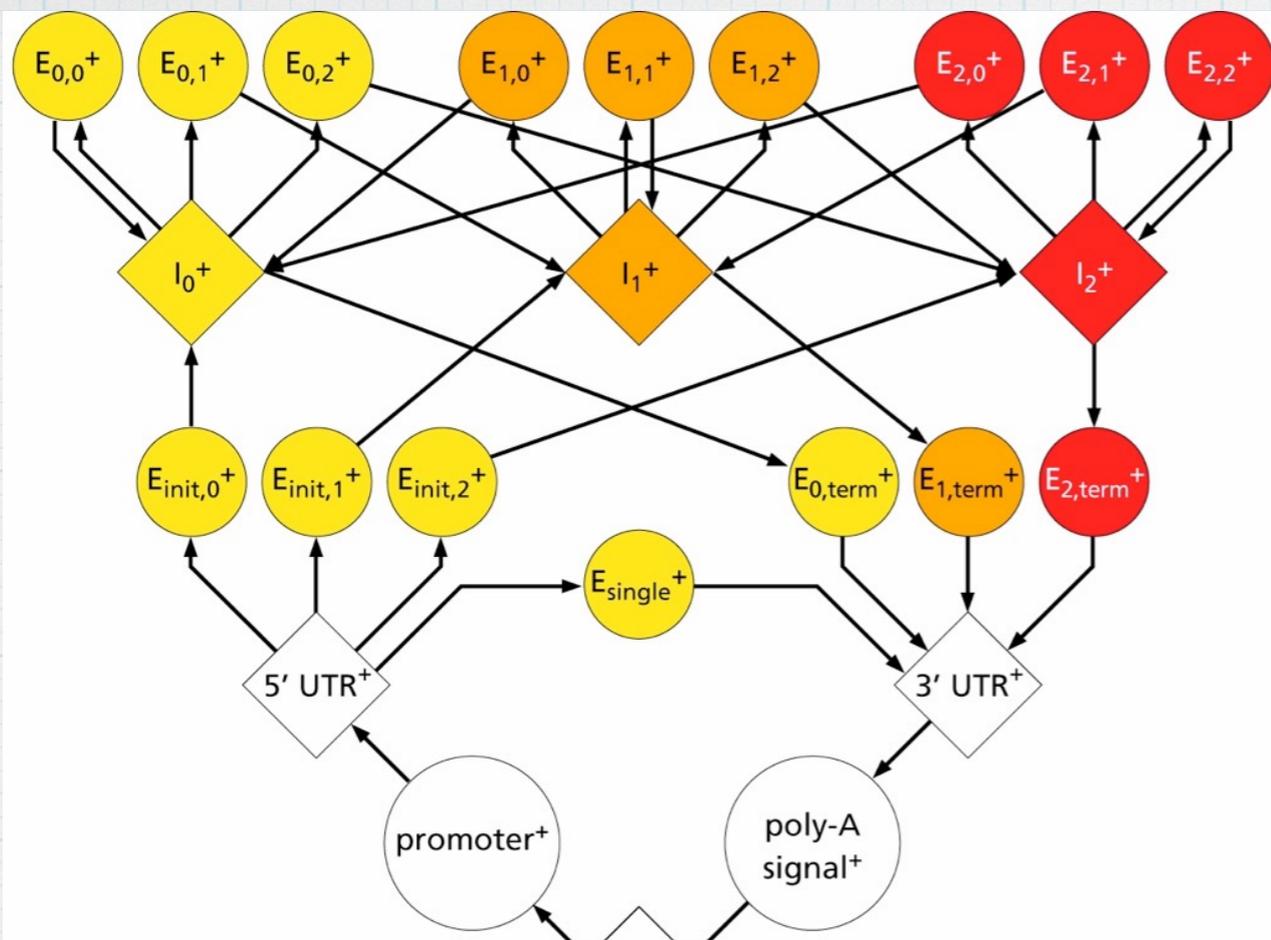
	sequence position									
	-3	-2	-1	+1	+2	+3	+4	+5	+6	
observed bases	A%	33	60	8	0	0	49	71	6	15
	C%	37	13	4	0	0	3	7	5	19
	G%	18	14	81	100	0	45	12	84	20
	U%	12	13	7	0	100	3	9	5	46
	consensus	A/C	A	G	G	U	A/G	A	G	U

(B)



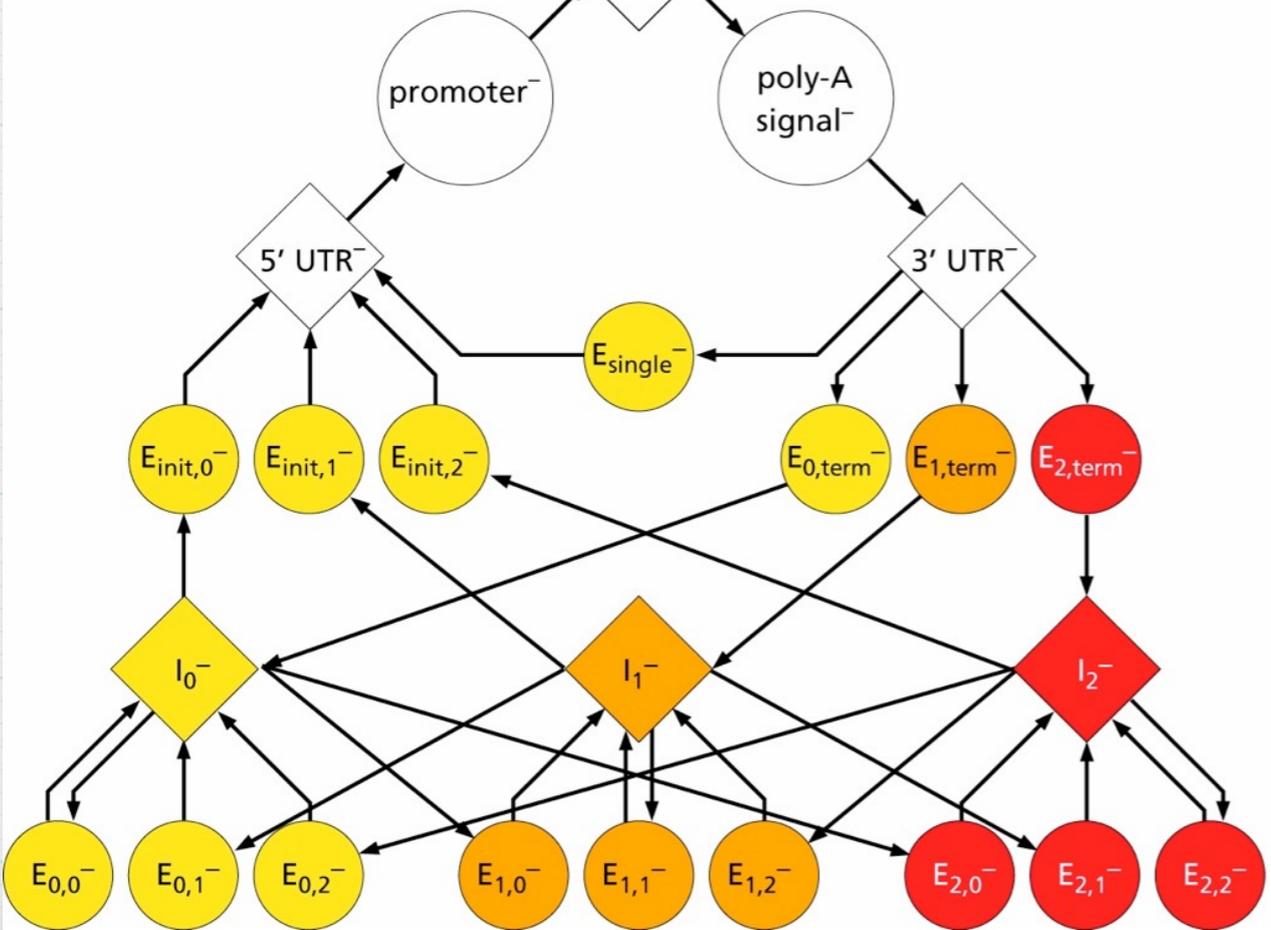
# Complete Eukaryotic Gene Models

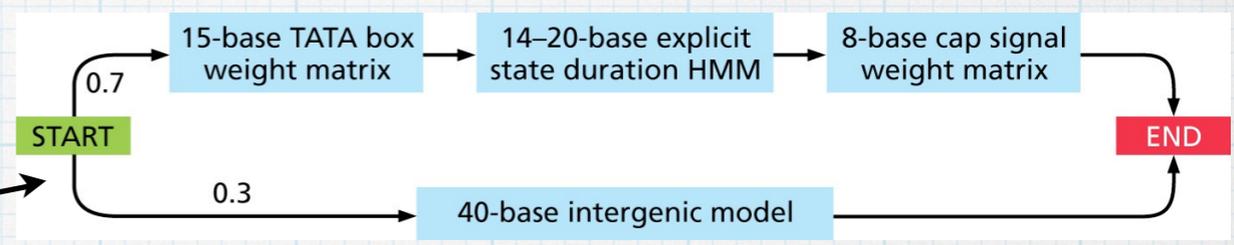
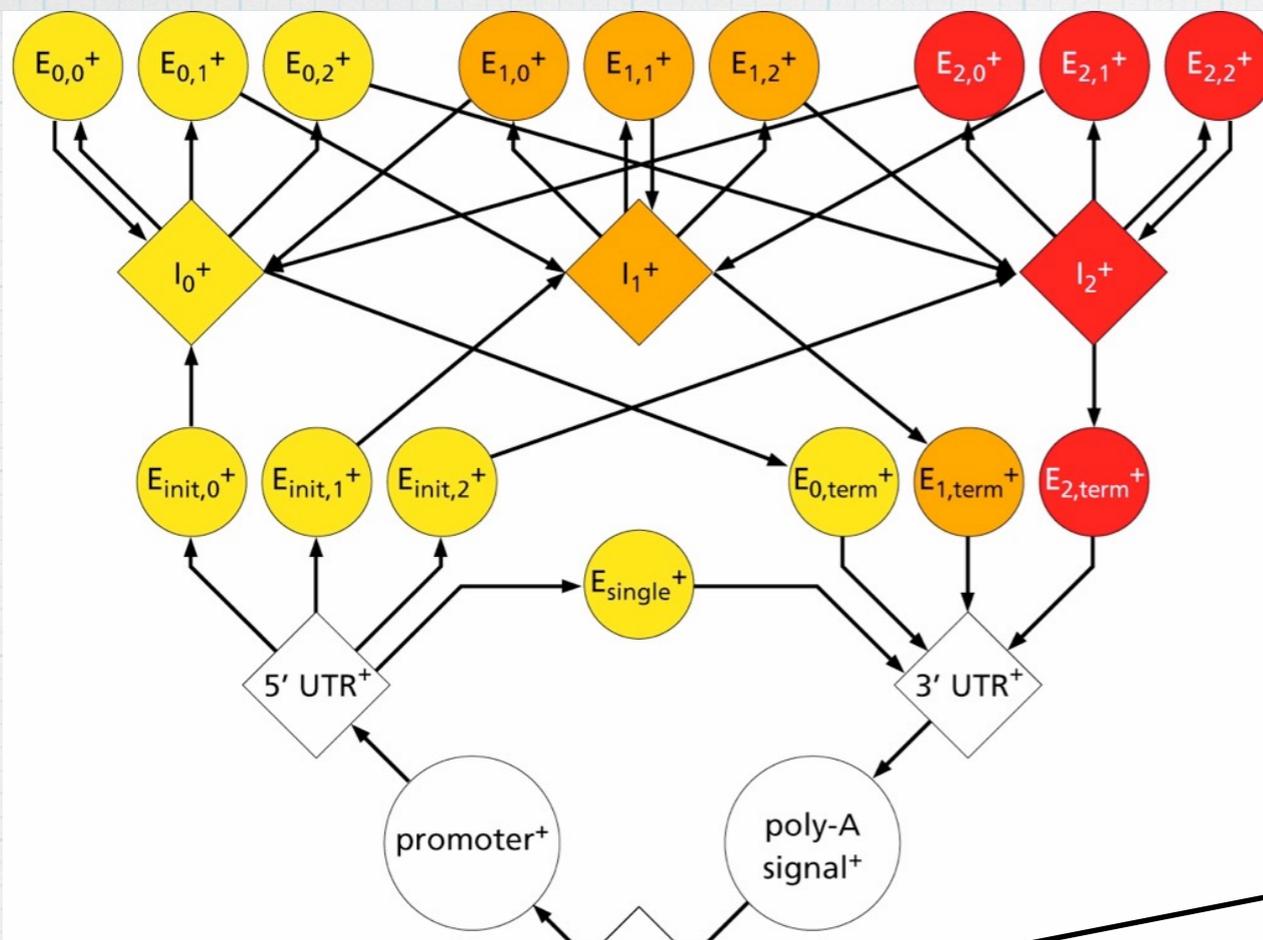
- \* Once all the separate components of a gene have been predicted, it is possible to put them all together to predict a complete gene structure.
- \* GenScan uses an HMM that considers both the forward and backward strands of DNA simultaneously.



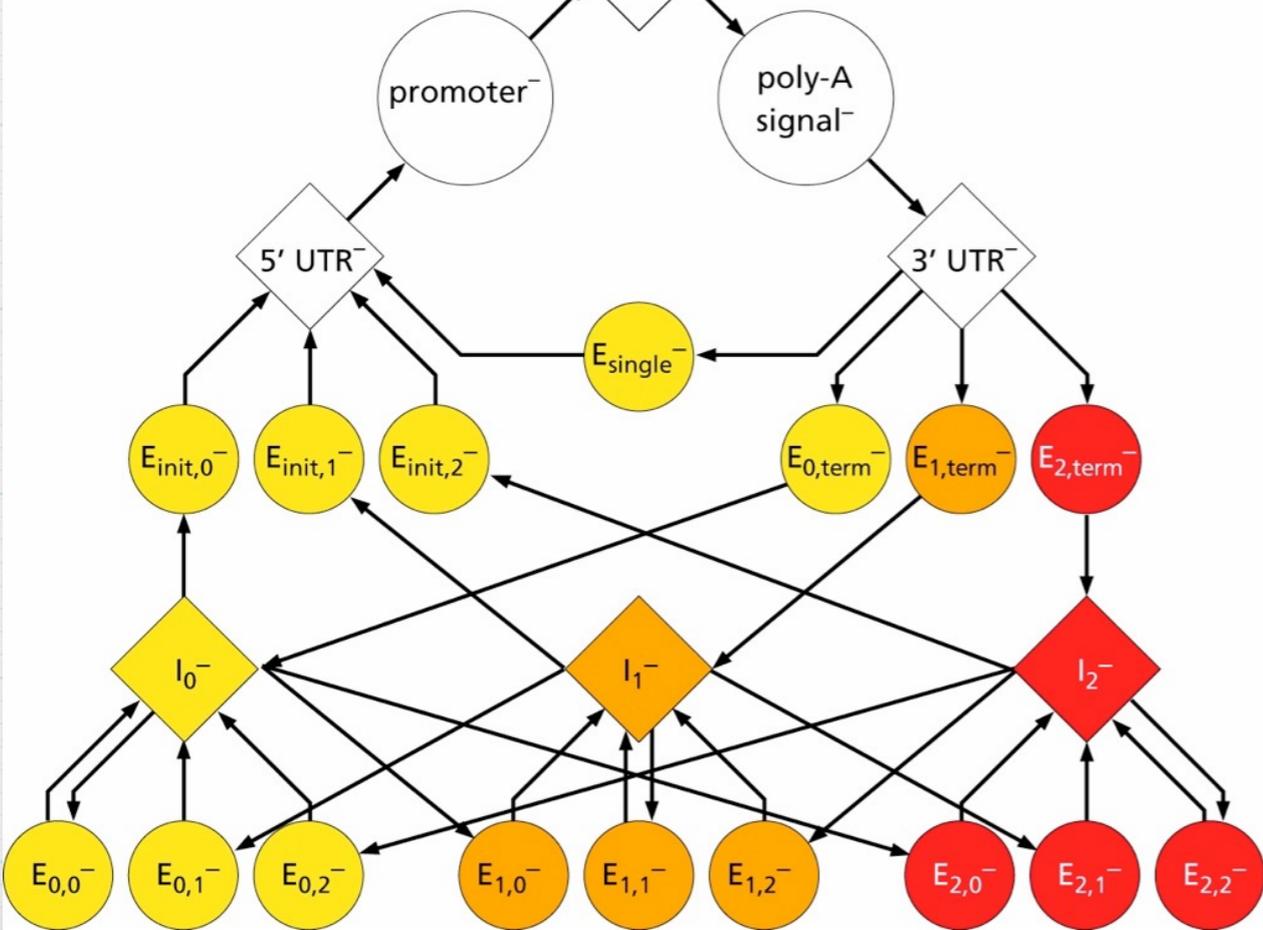
FORWARD (+) STRAND

REVERSE (-) STRAND





FORWARD (+) STRAND  
REVERSE (-) STRAND

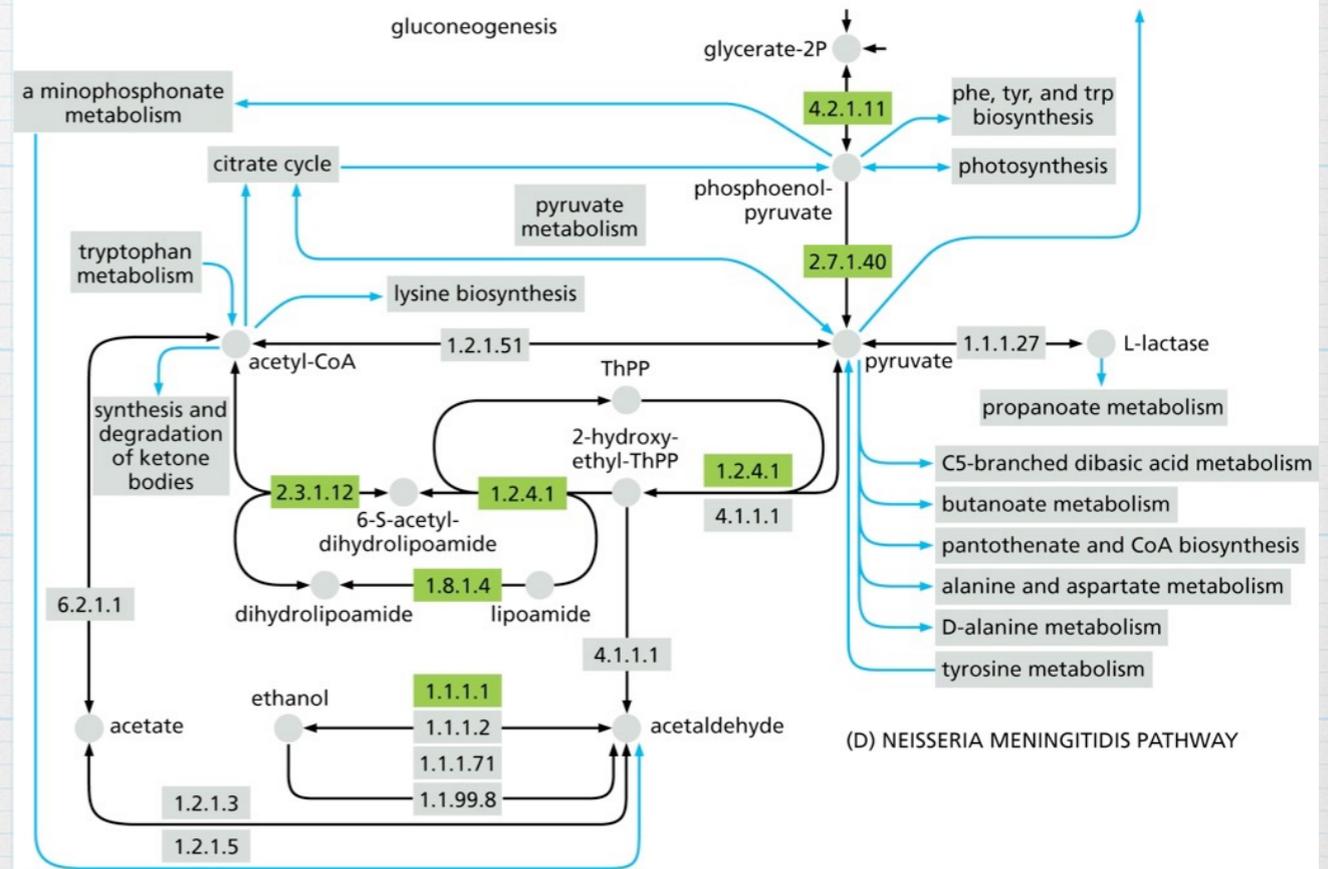
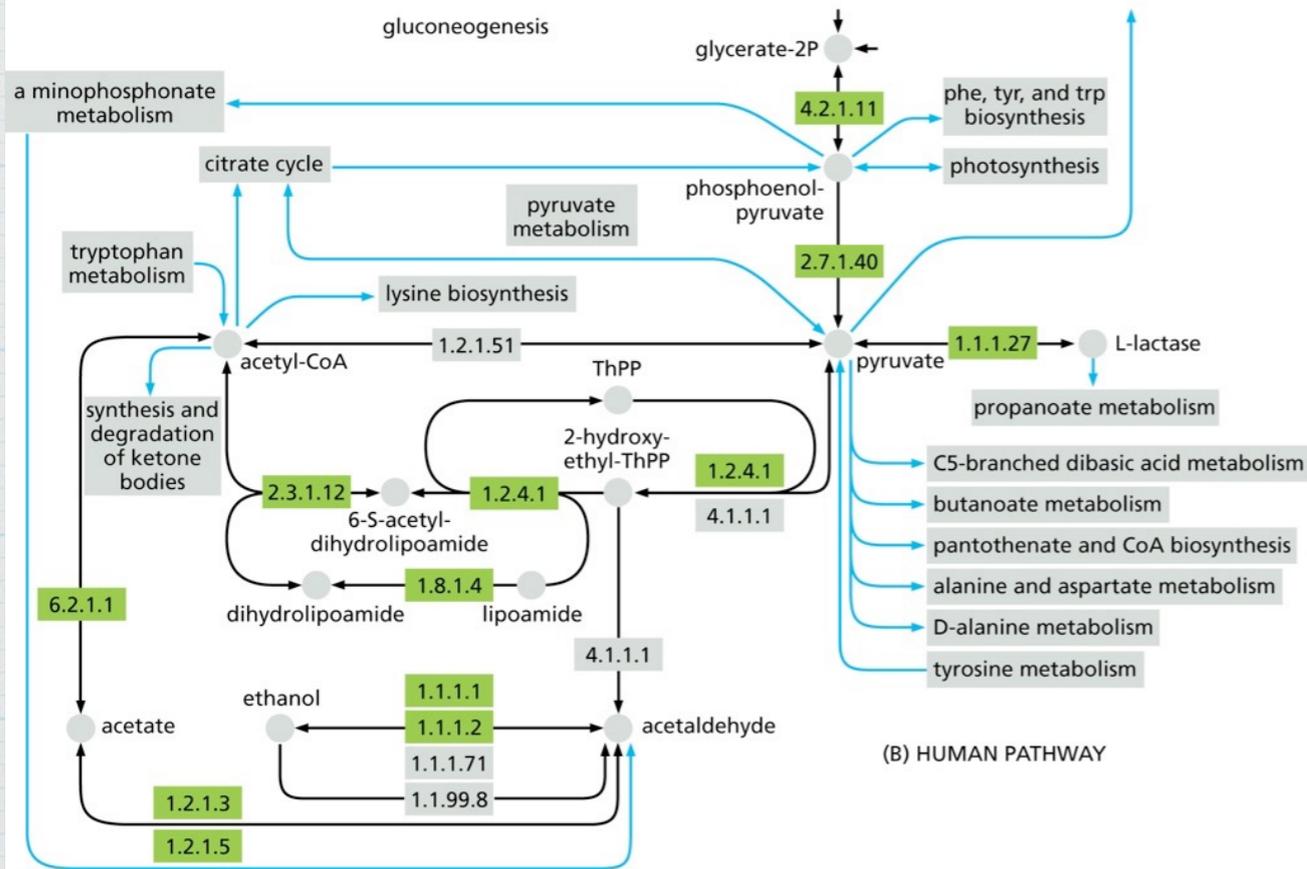
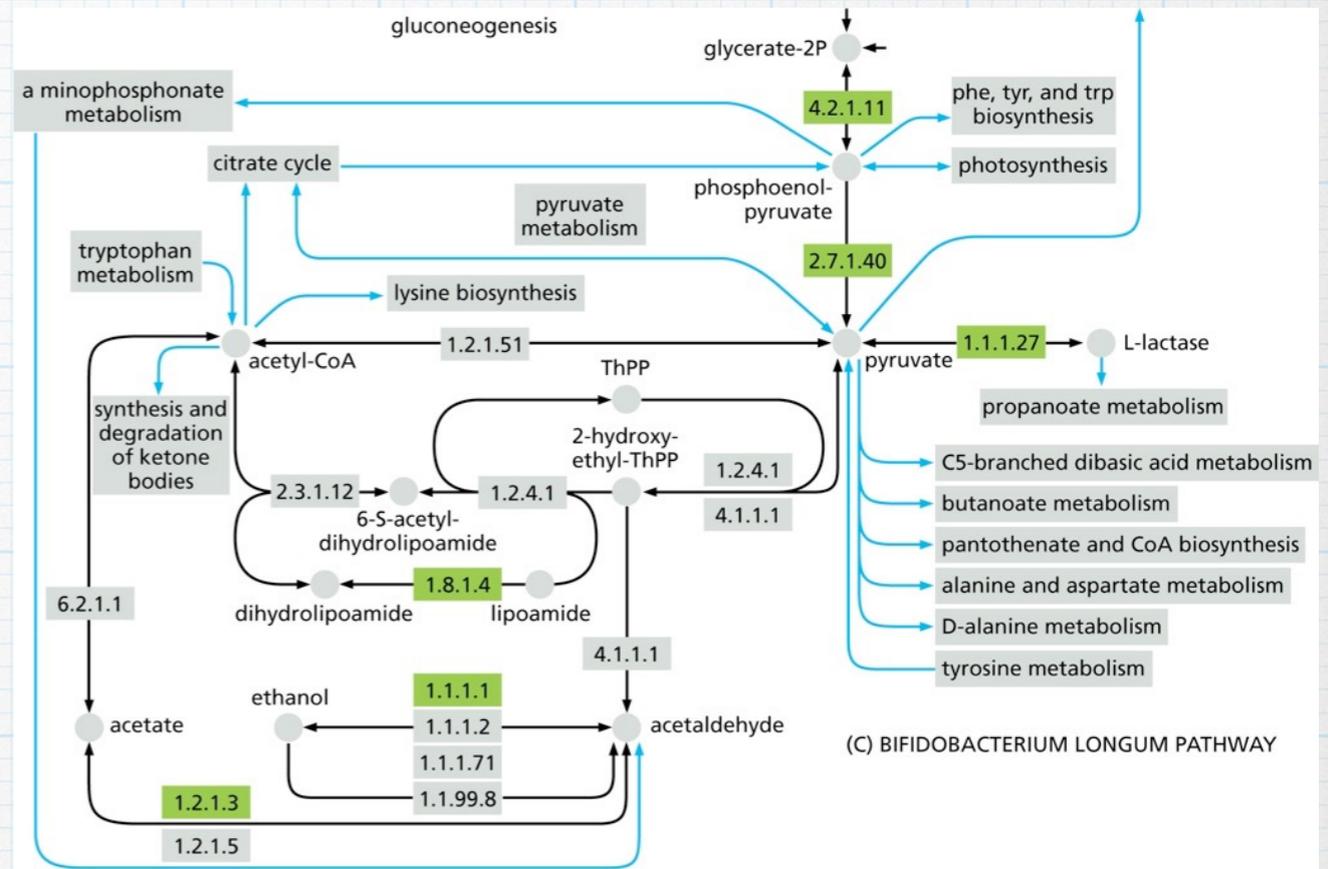
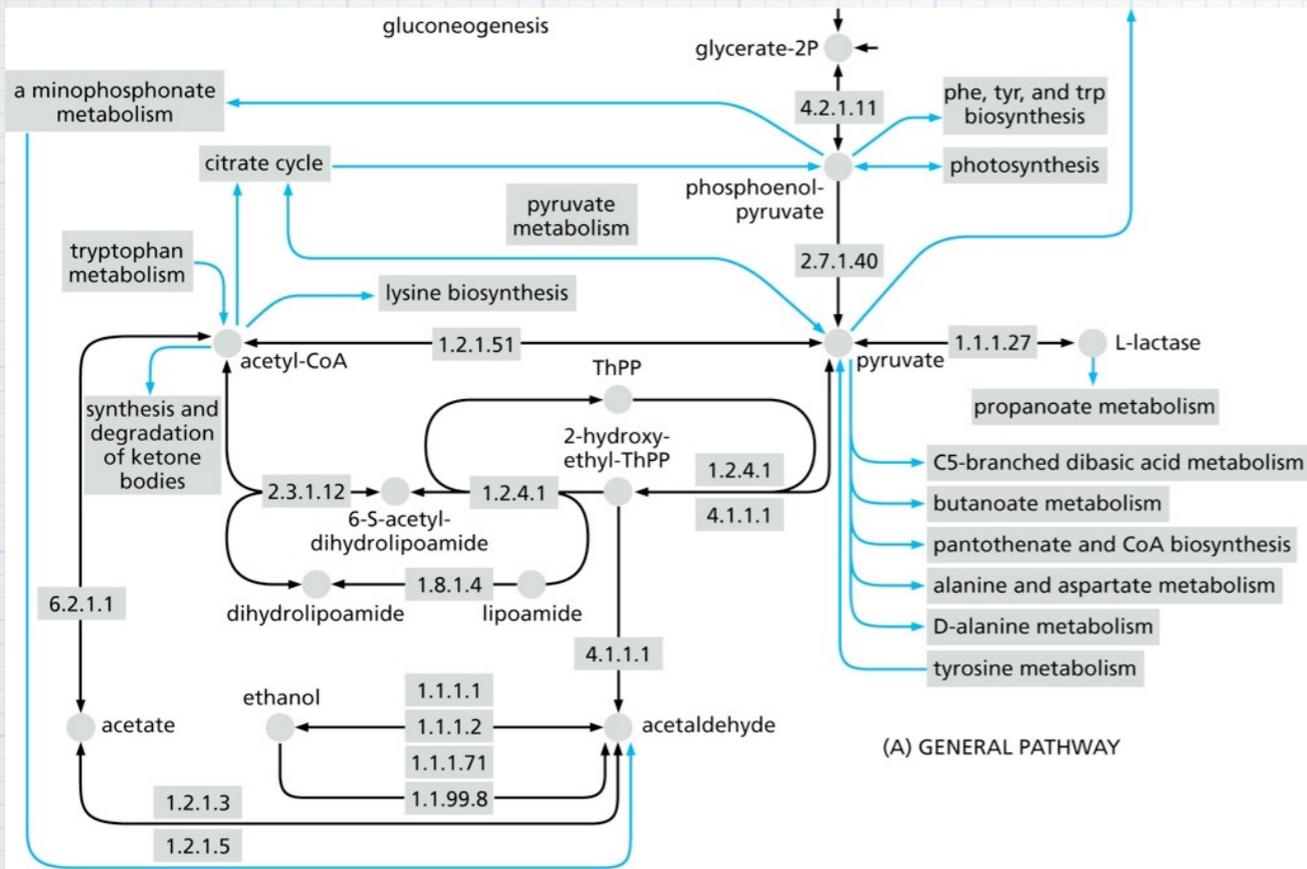


# Genome Annotation

- \* Once all the genes have been predicted, it remains to be determined what function the encoded proteins might play.
- \* The obvious way to start to determine gene function is by sequence analysis.
- \* If the encoded protein has one or more significant matches against sequence and pattern databases, the function and other properties can be predicted with considerable confidence to be similar to those of the matches.

# Pathway Information

- \* The use of pathway information can aid gene and genome annotation.
- \* Comparing a new genome to a well-annotated and functionally defined one can aid the analysis of specific pathways and may identify missing components or blanks in the pathway.



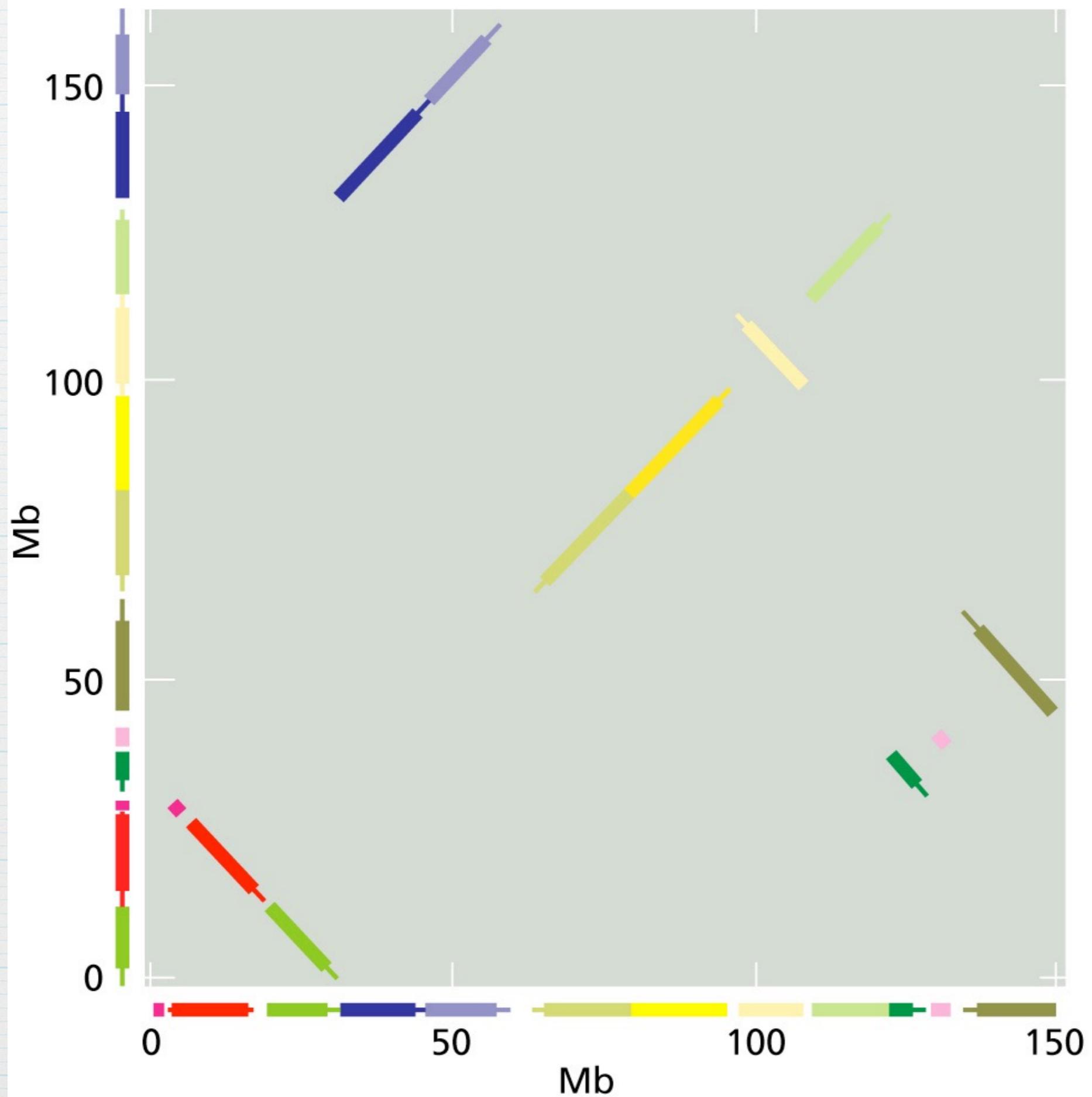
# Gene Ontology

- \* One of the important aspects of genome annotation has been the recognition of the importance of gene ontology.
- \* Gene ontology is a set of standardized and accepted terms that encompass the range of possible functions and can be found on the Gene Ontology Consortium's website ([geneontology.org](http://geneontology.org))

# Genome Comparison

- \* Comparison of two genomes can be a very powerful tool in determining the status of uncertain gene predictions.
- \* Aligning two genomes is not easy, as large-scale rearrangements are common, but it should be possible to find regions of synteny where the gene structure is sufficiently similar as to make their common evolutionary ancestry apparent.

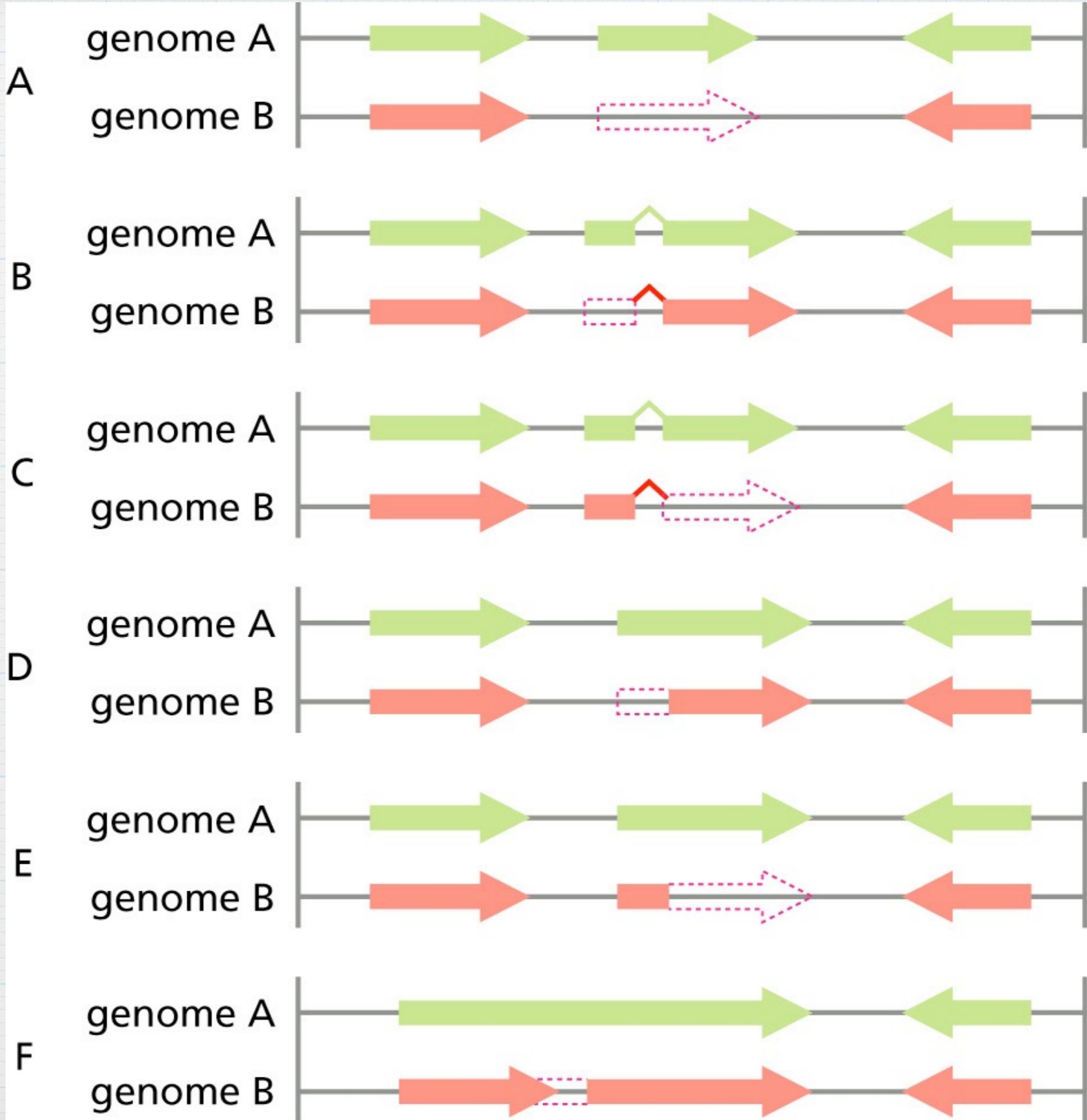
mouse chromosome X



rat chromosome X

# Genome Comparison

- \* Comparing syntenic regions can reveal errors in one or other genome annotation.



# Acknowledgments

- \* “Understanding Bioinformatics” by Zvelebil and Baum.