

MIT OpenCourseWare
<http://ocw.mit.edu>

6.047 / 6.878 Computational Biology: Genomes, Networks, Evolution
Fall 2008

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.

**Computational Gene Prediction and
Generalized Hidden Markov Models**

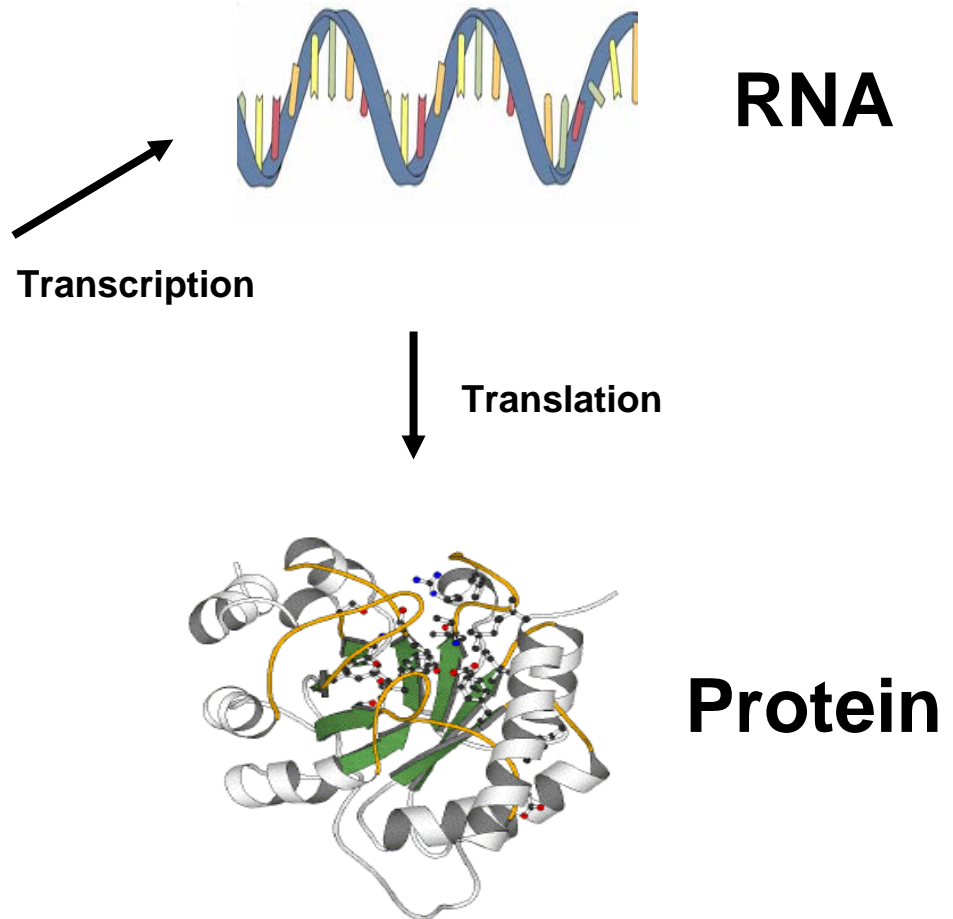
Today

- Gene Prediction Overview
- HMMs for Gene Prediction
- GHMMs for Gene Prediction
- Genscan

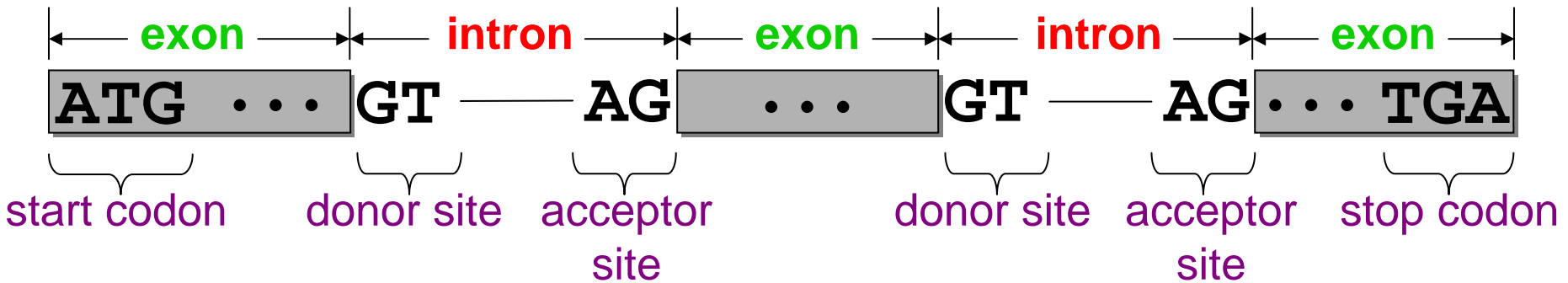
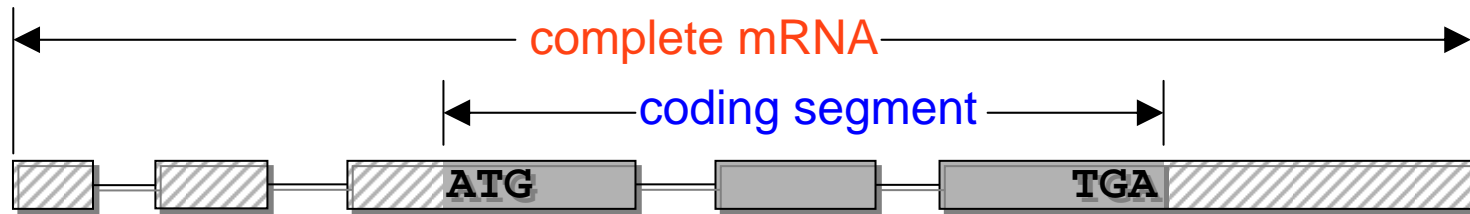
Genome Annotation

Genome Sequence

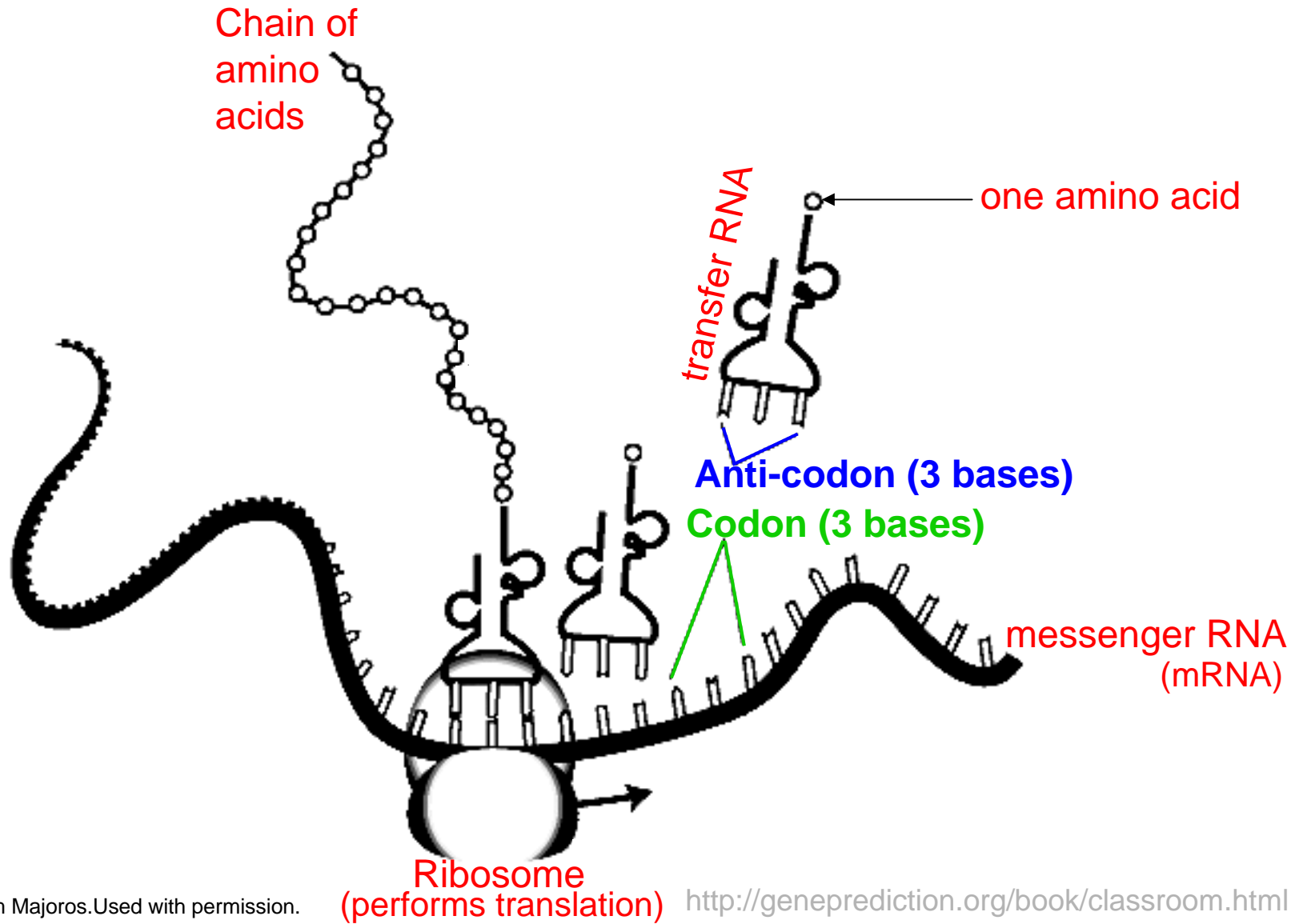
```
> HSKIIIBE, Human gene for casein kinase II subunit beta (EC 2.7.1.37).
ggggctgagatgtaaatagaggagctgggagggagtgcctcagagtttgggttgctttaagaaaggt
ggttccgaattctcccgtggttggagggccgaatgtgggaggaggaggataccagaggcaggaagga
gaactgagcttactgacactgttcttttctagctgacgtgaagatgagcagctcagaggaggtgc
ctggatttctctggtctctgtgggctccgtggcaatgaaattctctgtggaagtgccttcaacctcc
ctacttgccagcttcacatactctcccaccagcgttctctacacatttccacttctacactgttctct
aaagctttatgggagagagtgtaggtgaactaggagagacacaagtaactctgctgagttgggagtg
agaaacaagcacaacagatgcagttgtgtgtagatgaaggcatcacttagagcatttgcccagggtcaa
agatgaggattttgataggggttccctctggcttccatgtcctgacaggtggatgaagactacatcca
ggacaaatttaacttactggactcaatgagcaggtcccctcactatcgacaagctctagacatgatctt
ggacctggagcctggtgaggcaccctcagggttgtttgtgtgtgtgcgtgcactattttctctcaca
atctctattcacttgccctgaatttgccaaatttcttgggtctctgatttctttaaccccaaatca
tgctttatttgatcctccacctgactctgtctagtttgtgacgtatcactgttctcagtttt
tgcaagggtcagaagcccaggttctctgggtccatgcccagatgtggaagggttaaggcccaaaagta
ggtgctagccaaactgaatagcccgcagcccctggatagggcagggcaccctaggaagcgtgaaaaaa
agtgttgcatttggcccggctgtggttcagatgaagaactggaagacaaccccaaccagagtgcactg
attgagcaggcagccgagatgctttatggatgataccacgcccgtacatccttaccacccgtggcctc
gcccagatggtgagggctctctgctctacctgcctctctgagcagtaagagacacaggttctctgca
gcaagaagtcatgtttaagccctgtttaaggaagctagctgagaagaggggaagaaaccccaagacttg
ccctgcccctaaattggaagaaggaacacacagaagttagagagcccactagtcagagagaagggggcct
ctggacagagtggaaaggagtggcagacagagtgggtatgggttgggctgccaagggagttgcccctct
ttacatctacctgccaaccccttccatgttatcaccctcagttggaagaactaccagcaaggagacttgg
gttactgtctcgtgtgtactgtgagaaccagccaatgcttcccatGGtgagtggtgagaaggggaaa
ggaaagcaccgttggcagctctatgggaaggagtggggctcaacacatggagcctgagctcctgagg
ggaggttaggttaggaataggggataacctggcctgctgagctggctgtctcccaggccttccagacat
cccaggtgaaagccatggtgagctctactgccccaaagtgcattggtgtacacacccaagtcatacaag
acaccatcacacaggatggcgccctacttcggcactggtttccctcactgctctcactggtgcatccca
gtacccggcccaagagacctgccaaccagtttgtgcccaggtagggagcagggagagtcattaaaggtca
aaggaaagcccagaatcccccaagagaggggaggaagggcagggcattggcccttcttgaggtctgctctcc
cagaaatcagggcatctcccctgctgagtgactgtgggaaagtatttgatattctgctgagttacct
tatgtagaatgtcttgagctgagaagtgggaaccacagggcttagctctgagcaggtccaatagag
gagctcaggtggggaggtgggaaatgcaggtgactggcagggcctggaaggggctcaatgctgctgcccct
ctgacctctgcccctggcctaggctctacgggttccaagatccaatccgaaggcctaccagctgacgctcca
agccggcagcaactccaagagccagtcagaagcagatcgctgatccctcccccaactgctcctgagctc
tttctcttttctctcttttggccaccttccaggaacccctgatggttttagtttaaataaagga
gtcgttatcgtggtgggaaatgaaaataaagtagaagaaaaggccaagcagctgctgctggtgcttgc
bgaaggggggtgagcgtggccatggaaatcgggctccaagggccagggtatgg
```



Eukaryotic Gene Structure



Translation



Genetic Code

Acid	Codons	Acid	Codons	Acid	Codons	Acid	Codons
A	GCA GCC GCG GCT	G	GGA GGC GGG GGT	M	ATG	S	AGC AGT TCA TCC TCG TCT
C	TGC TGT	H	CAC CAT	N	AAC AAT	T	ACA ACC ACG ACT
D	GAC GAT	I	ATA ATC ATT	P	CCA CCC CCG CCT	V	GTA GTC GTG GTT
E	GAA GAG	K	AAA AAG	Q	CAA CAG	W	TGG
F	TTC TTT	L	CTA CTC CTG CTT TTA	R	AGA AGG CGA CGC CGG CGT	Y	TAC TAT

Each amino acid is encoded by one or more *codons*.

Each codon encodes a single *amino acid*.

The *third position* of the codon is the most likely to vary, for a given amino acid.

Figure by MIT OpenCourseWare.

Gene Prediction as “Parsing”

- Given a genome sequence, we wish to label each nucleotide according to the parts of genes
 - Exon, intron, intergenic, etc
- The sequence of labels must follow the syntax of genes
 - e.g. exons must be followed by introns or intergenic not by other exons
- We wish to find the optimal parsing of a sequence by some measure

Features

A *feature* is *any DNA subsequence of biological significance*.

For practical reasons, we recognize two broad classes of features:

signals — short, fixed-length features

content regions — variable-length features

Content Regions

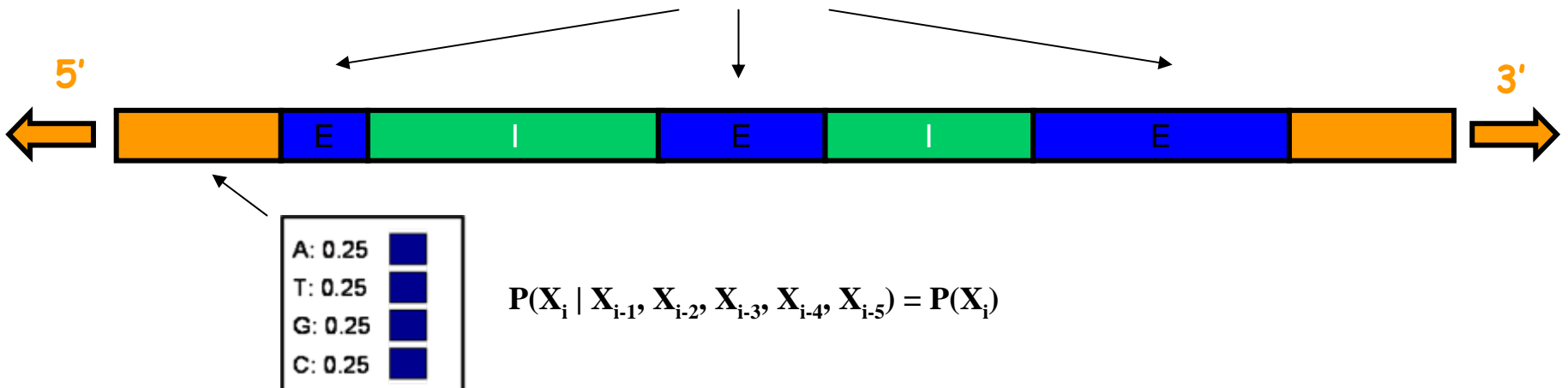
Content regions often have characteristic base composition

Example

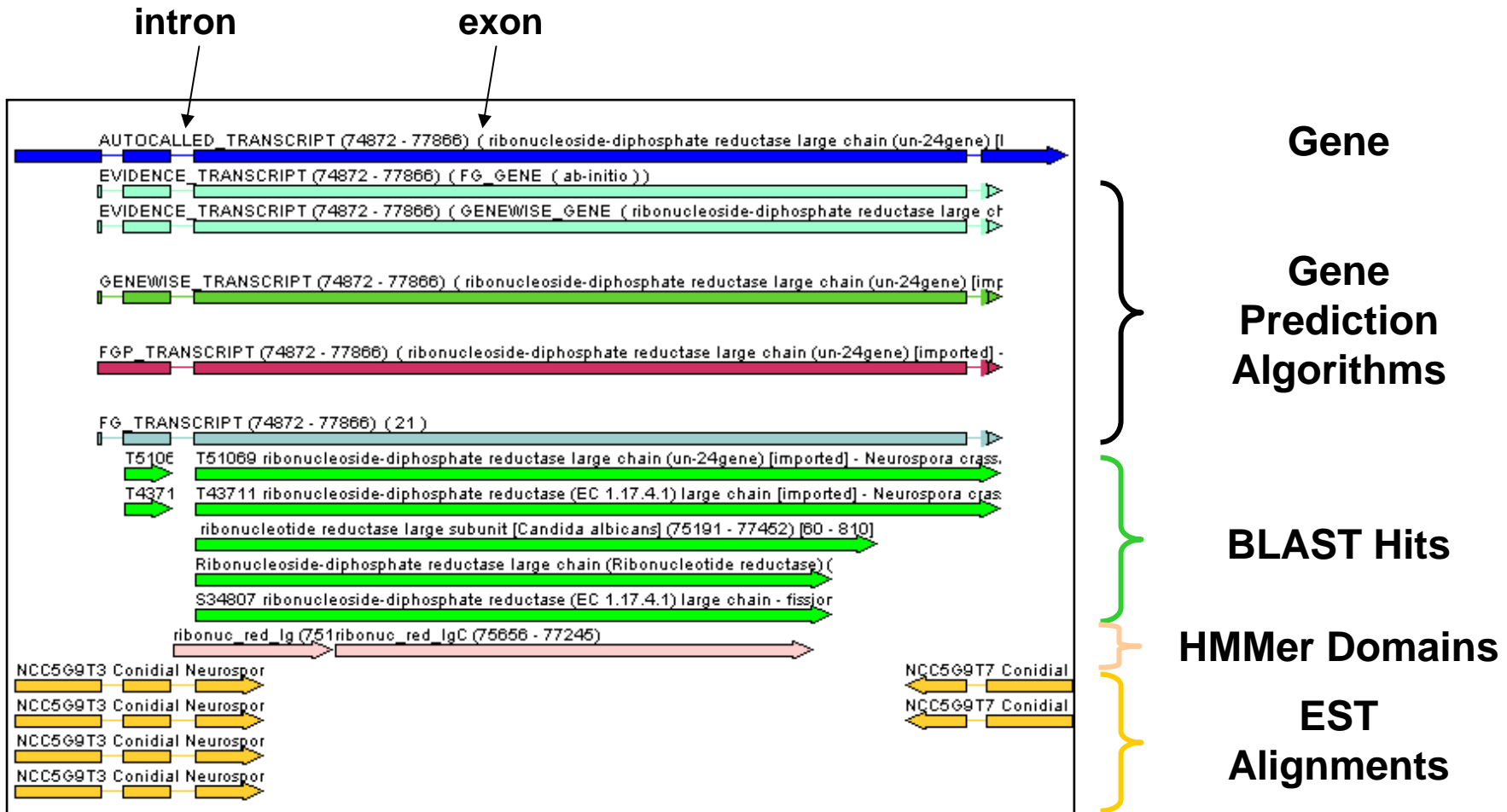
- Recall: often multiple codons for each amino acid
- All codons are not used equally

Characteristic higher order nucleotide statistics in coding sequences
(hexanucleotides)

$$P_{\text{exon}}(X_i | X_{i-1}, X_{i-2}, X_{i-3}, X_{i-4}, X_{i-5})$$



Extrinsic Evidence



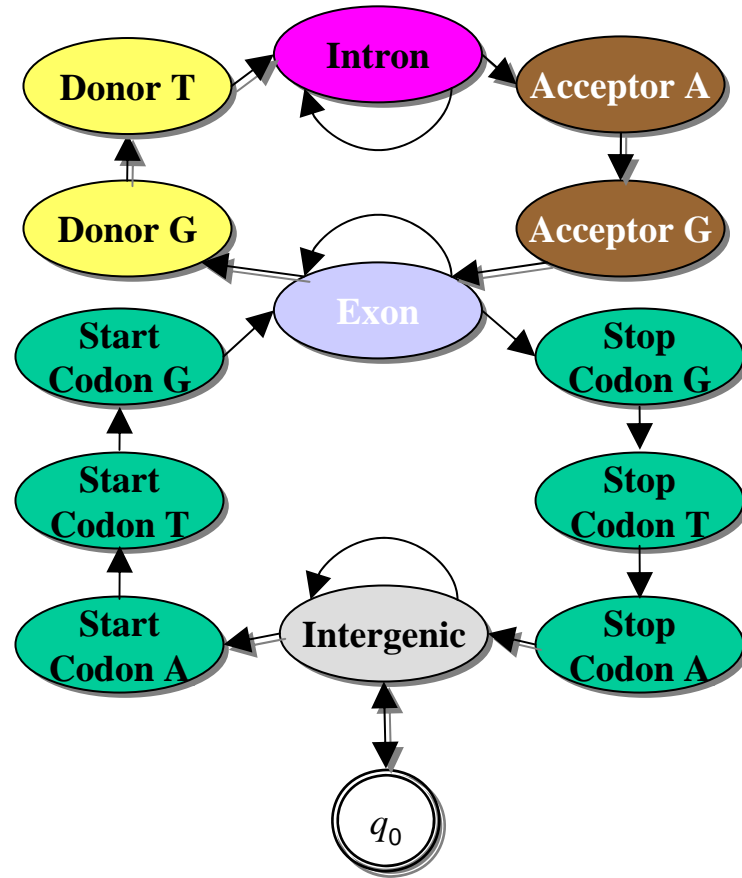
Neurospora crassa (a fungus)

HMMs for Gene Prediction

- States correspond to gene and genomic regions (exons, introns, intergenic, etc)
- State transitions ensure legal parses
- Emission matrices describe nucleotide statistics for each state

A (Very) Simple HMM

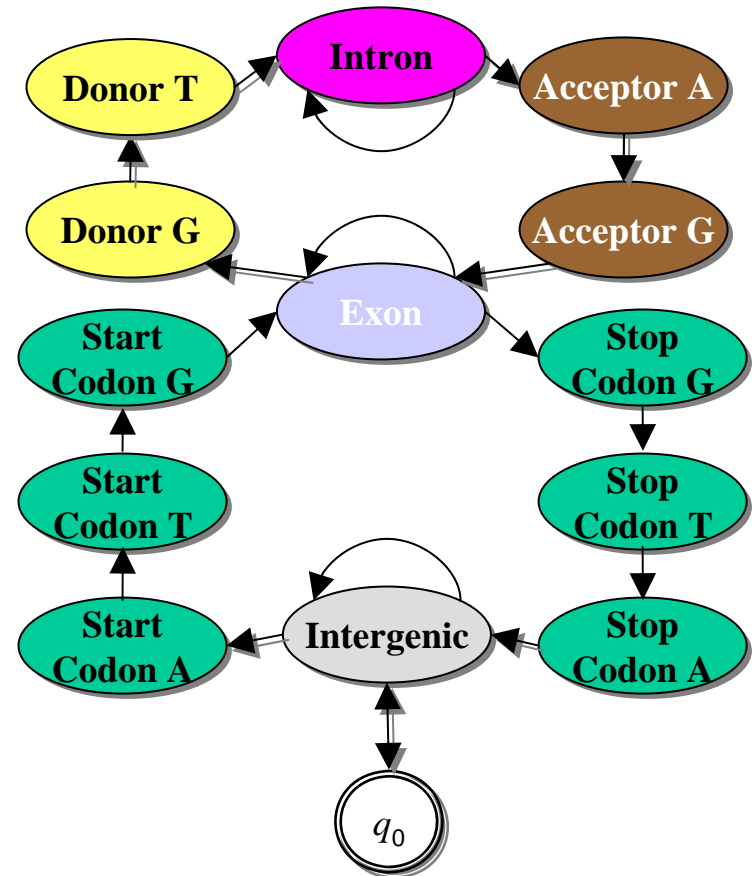
the Markov model:



A Generative Model

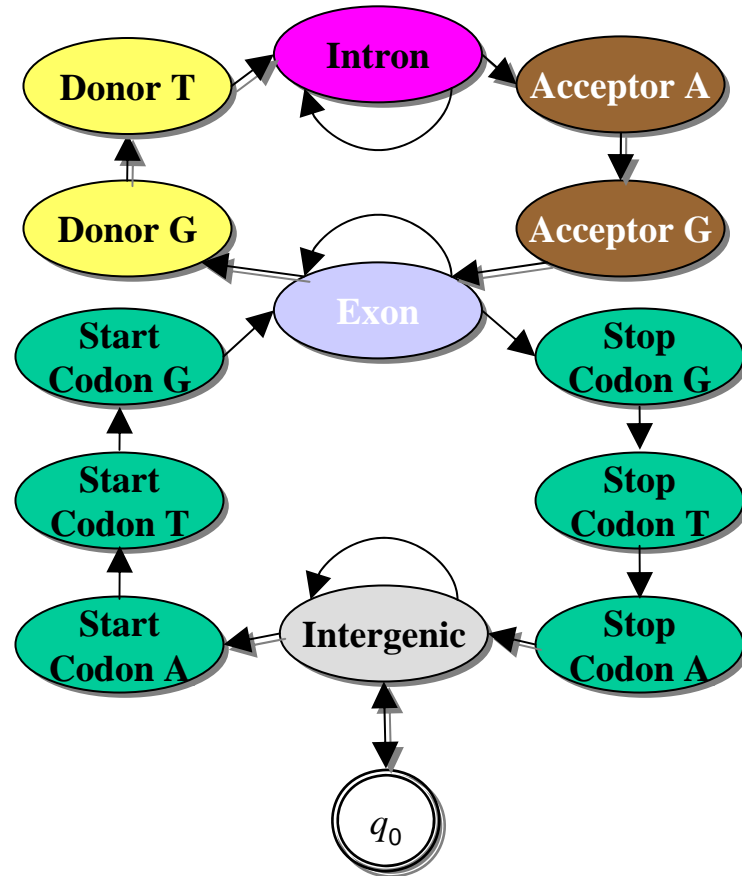
We can use this HMM to generate a sequence and state labeling

- The initial state is q_0
- Choose a subsequent state, conditioned on the current state, according to $a_{jk} = P(q_k | q_j)$
- Choose a nucleotide to emit from the state emissions matrix $e_k(X_i)$
- Repeat until number of nucleotides equals desired length of sequence



But We Usually Have the Sequence

the Markov model:



the input sequence: AGCTAGCAGTATGTCATGGCATGTTCCGAGGTAGTACGTAGAGGTAGCTAGTATAGGTCGATAGTACGCGA

What is the best state labeling?

Finding The Most Likely Path

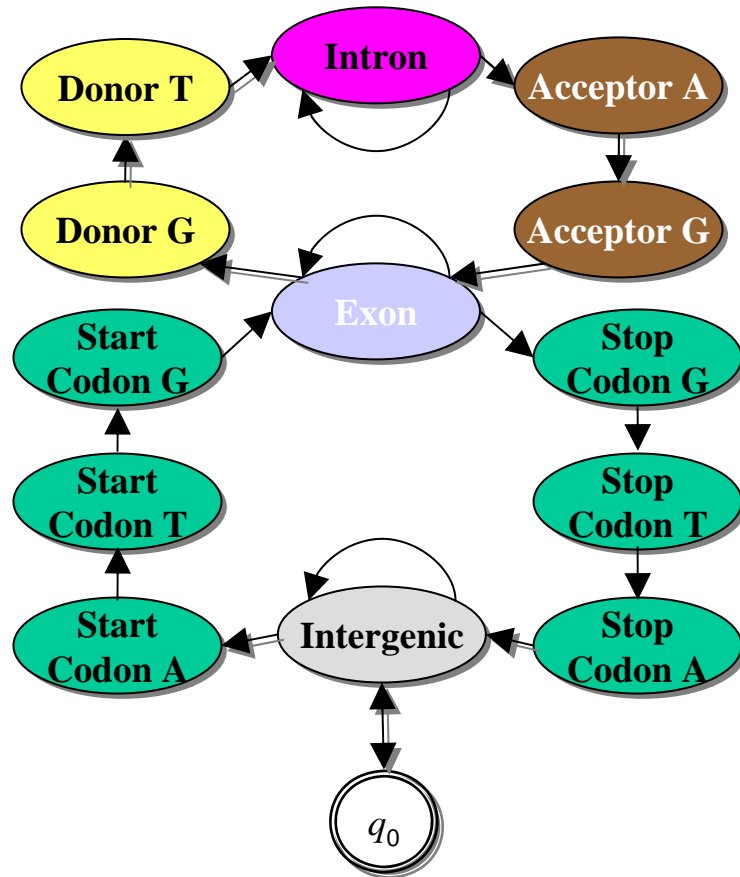
- A sensible choice is to choose π that maximizes $P[\pi|x]$
- This is equivalent to finding path π^* that maximizes total joint probability $P[x, \pi]$:

$$P(x, \pi) = \underbrace{a_{0\pi_1}}_{\text{start}}^* \prod_i \underbrace{e_{\pi_i}(x_i)}_{\text{emission}} \times \underbrace{a_{\pi_i\pi_{i+1}}}_{\text{transition}}$$

How do we select π^* efficiently?

A (Very) Simple HMM

the Markov model:



the input sequence: AGCTAGCAGTATGTCATGGCATGTTCCGAGGTAGTACGTAGAGGTAGCTAGTATAGGTCGATAGTACGCGA

the most probable path:

the gene prediction:

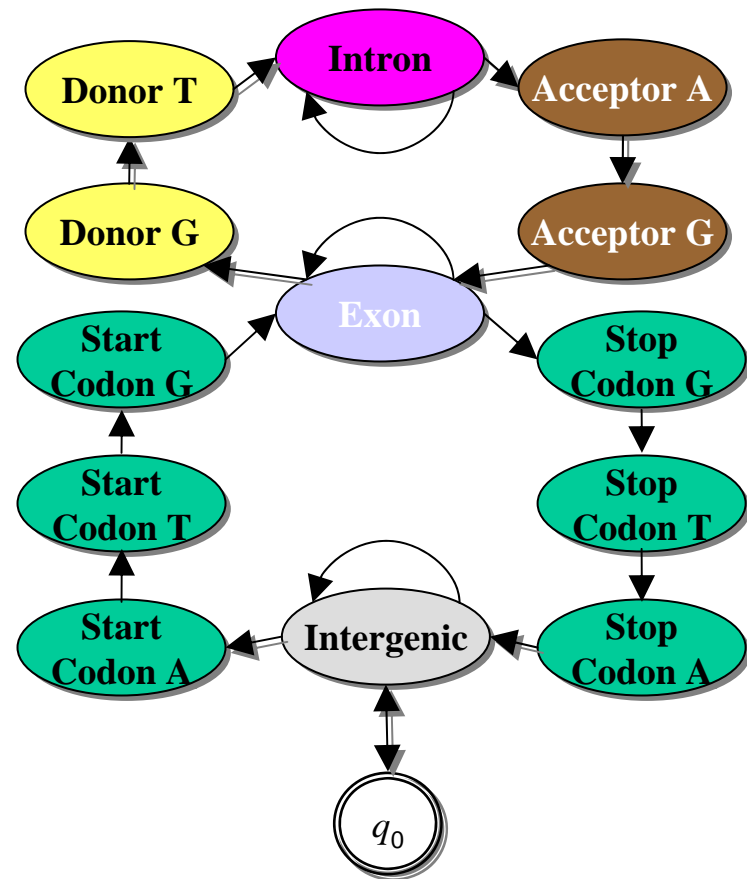
A Real HMM Gene Predictor

Title page of journal article removed due to copyright restrictions.
The article is the following:

Krogh, Anders, I. Saira Mian, and David Haussler. "A Hidden Markov Model That Finds Genes in *E.coli* DNA." *Nucleic Acids Research* 22, no. 22 (1994): 4768-4778.

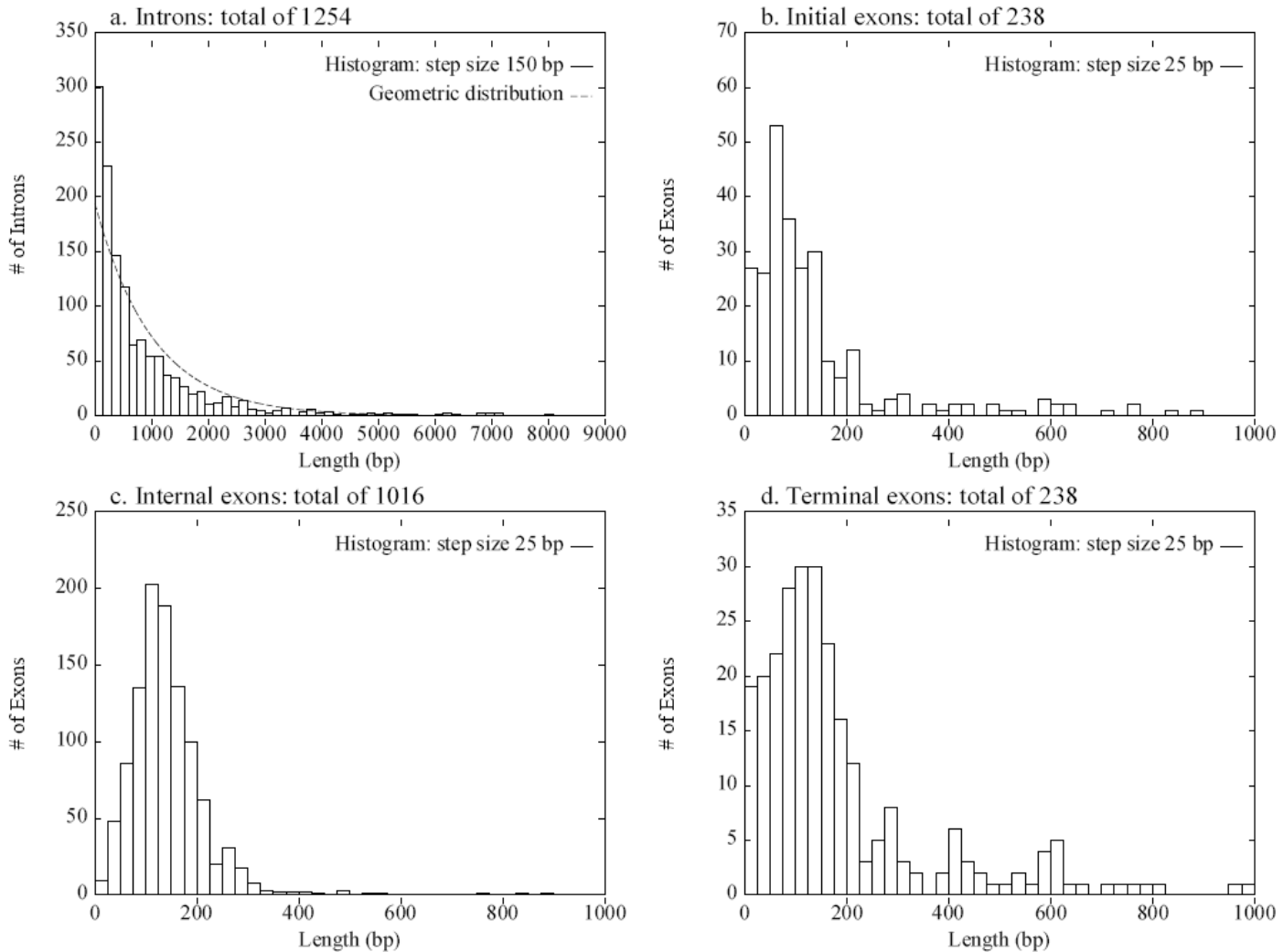
HMM Limitations

The HMM framework imposes constraints on state paths...



Human Exon Lengths

Fig. 1. Length distributions of introns and exons in human genes



Legend. Intron, exon length data from 238 multi-exon genes of GENSCAN learning set (Appendix A).

Courtesy of Christopher Burge. Used with permission.

Fungal Intron Lengths

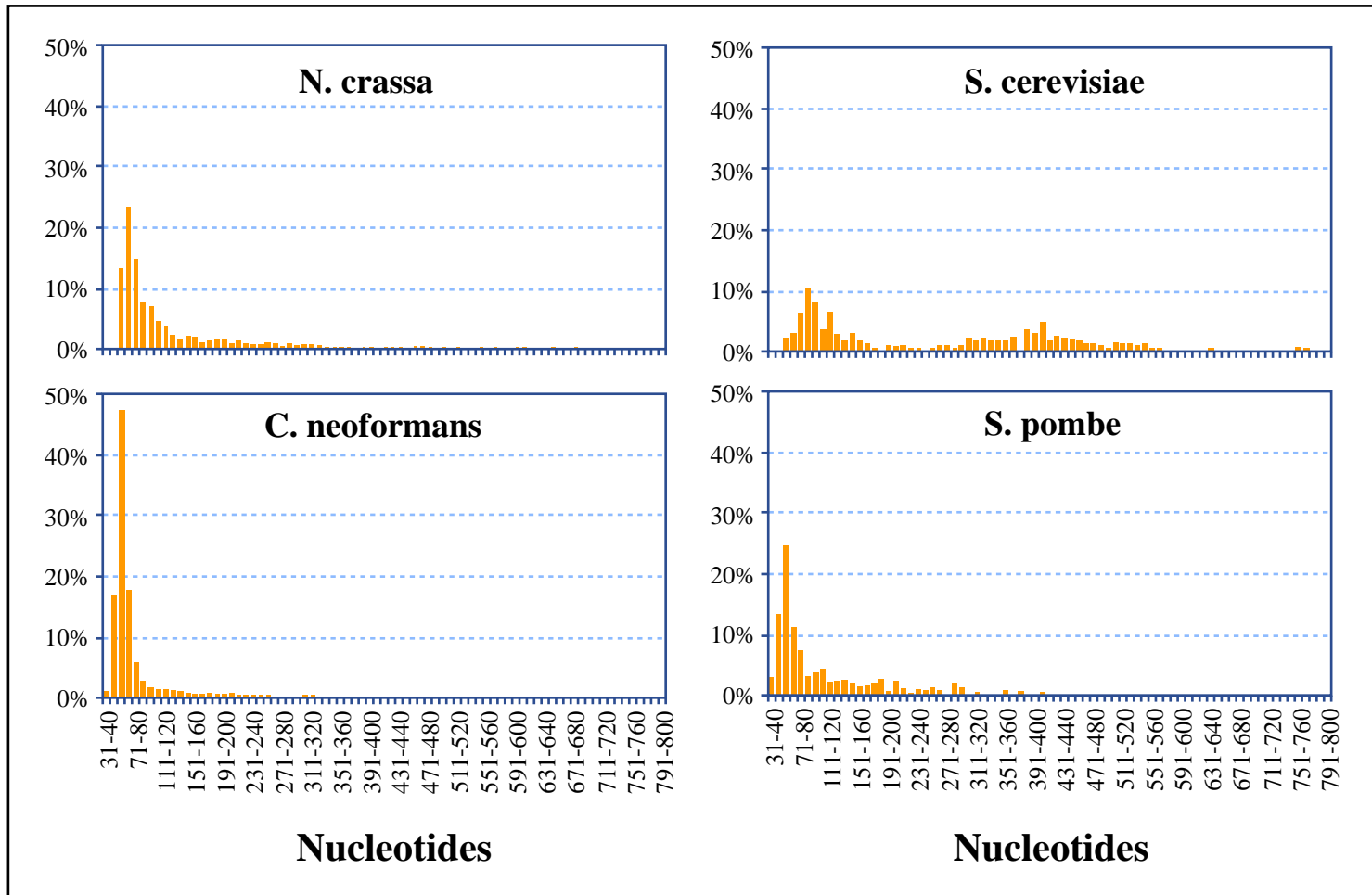
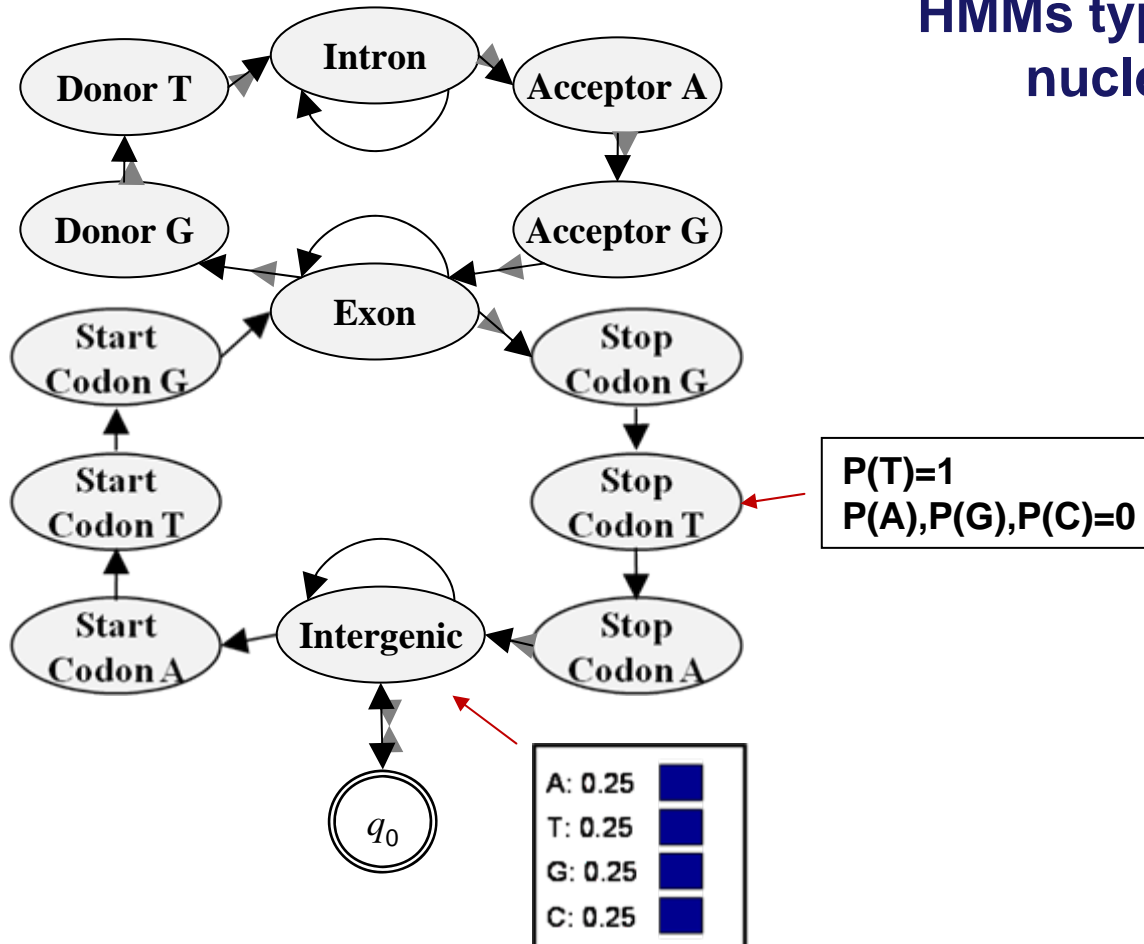


Figure by MIT OpenCourseWare.

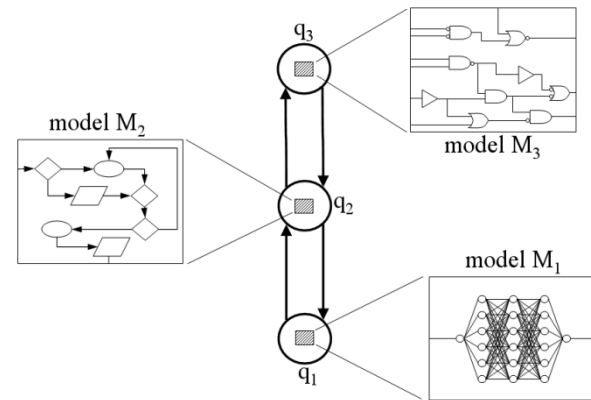
HMM Emissions

HMMs typically emit a single nucleotide per state



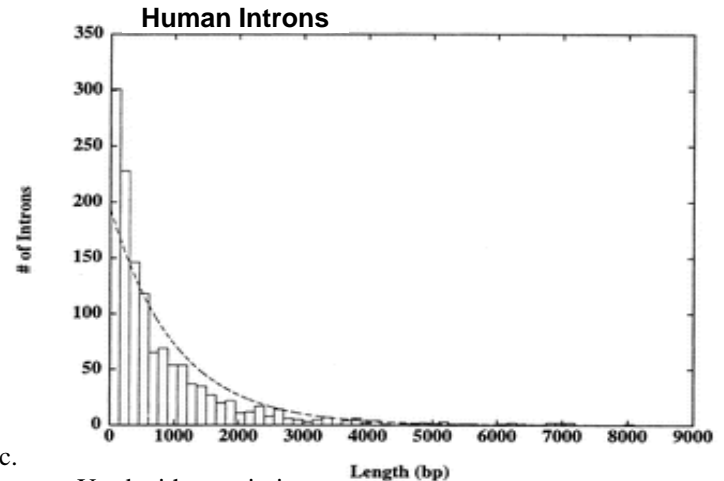
Generalized HMMs (GHMMs)

- GHMMs emit more than one symbol per state
- Emissions probabilities modeled by any arbitrary probabilistic model
- Feature lengths are explicitly modeled



W.H. Majoros (<http://geneprediction.org/book/classroom.html>)

Courtesy of William Majoros. Used with permission.



Courtesy of Elsevier, Inc.

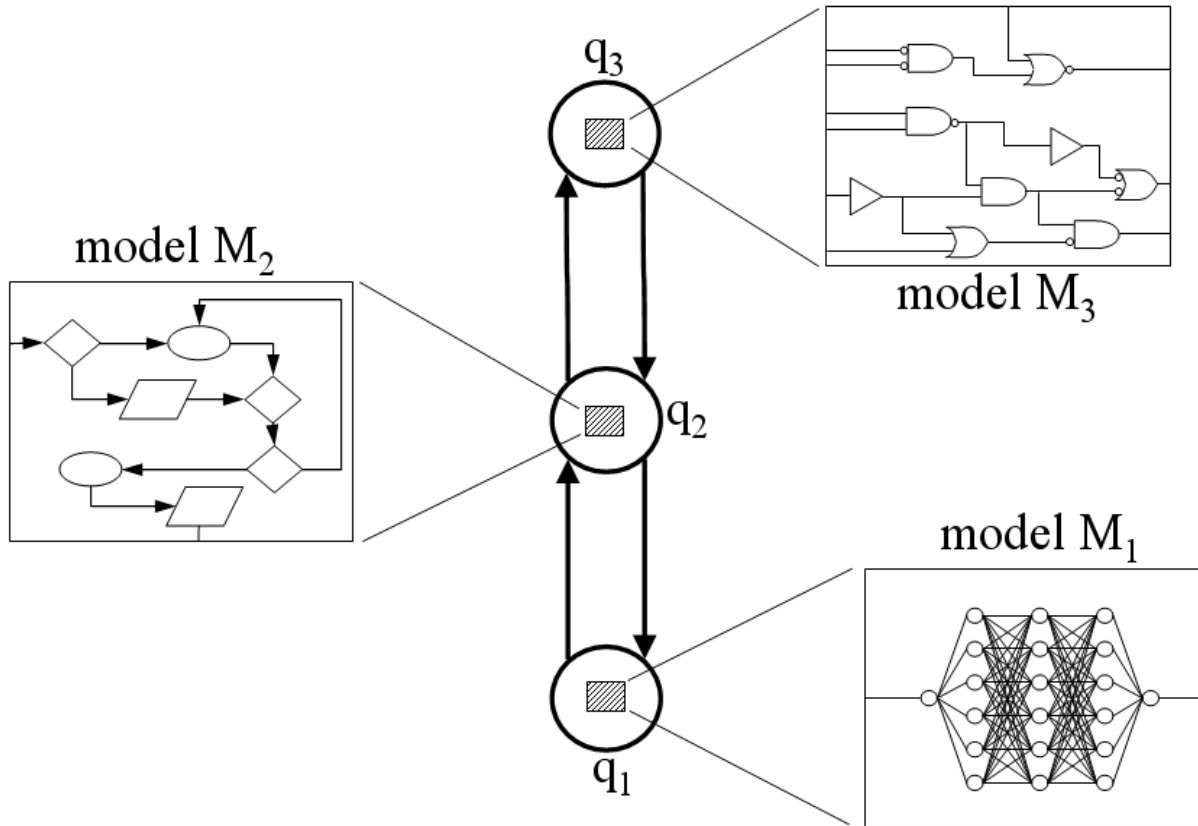
<http://www.sciencedirect.com>. Used with permission.

Burge, Karlin (1997)

GHMM Elements

- **States** Q
 - **Observations** V
 - **Initial state probabilities** $\pi_i = P(q_0=i)$
 - **Transition Probabilities** $a_{jk} = P(q_i=k|q_{i-1}=j)$
- } Like HMMs
- **Duration Probabilities** $f_k(d) = P(\text{state } k \text{ of length } d)$
 - **Emission Probabilities** $e_k(X_{\alpha, \alpha+d}) = P_k(X_{\alpha} \dots X_{\alpha+d} | q_k, d)$
↑
Now emit a subsequence

Model Abstraction in GHMMs



W.H. Majoros (<http://geneprediction.org/book/classroom.html>)

Models must return the probability of a subsequence given a state and duration

Courtesy of William Majoros. Used with permission.

GHMM Submodel Examples

1. **WMM (Weight Matrix)** $\prod_{i=0}^{L-1} P_i(x_i)$

2. **Nth-order Markov Chain (MC)** $\prod_{i=0}^{n-1} P(x_i | x_0 \dots x_{i-1}) \prod_{i=n}^{L-1} P(x_i | x_{i-n} \dots x_{i-1})$

3. **Three-Periodic Markov Chain (3PMC)** $\prod_{i=0}^{L-1} P_{(f+i) \pmod{3}}(x_i)$

5. **Codon Bias** $\prod_{i=0}^{n-1} P(x_{\alpha+3i} x_{\alpha+3i+1} x_{\alpha+3i+2})$

6. **MDD**

Ref: Burge C (1997) Identification of complete gene structures in human genomic DNA. *PhD thesis*. Stanford University.

7. **Interpolated Markov Model**

Ref: Salzberg SL, Delcher AL, Kasif S, White O (1998) Microbial gene identification using interpolated Markov models. *Nucleic Acids Research* 26:544-548.

$$P_e^{IMM}(s | g_0 \dots g_{k-1}) = \begin{cases} \lambda_k^G P_e(s | g_0 \dots g_{k-1}) + (1 - \lambda_k^G) P_e^{IMM}(s | g_1 \dots g_{k-1}) & \text{if } k > 0 \\ P_e(s) & \text{if } k = 0 \end{cases}$$

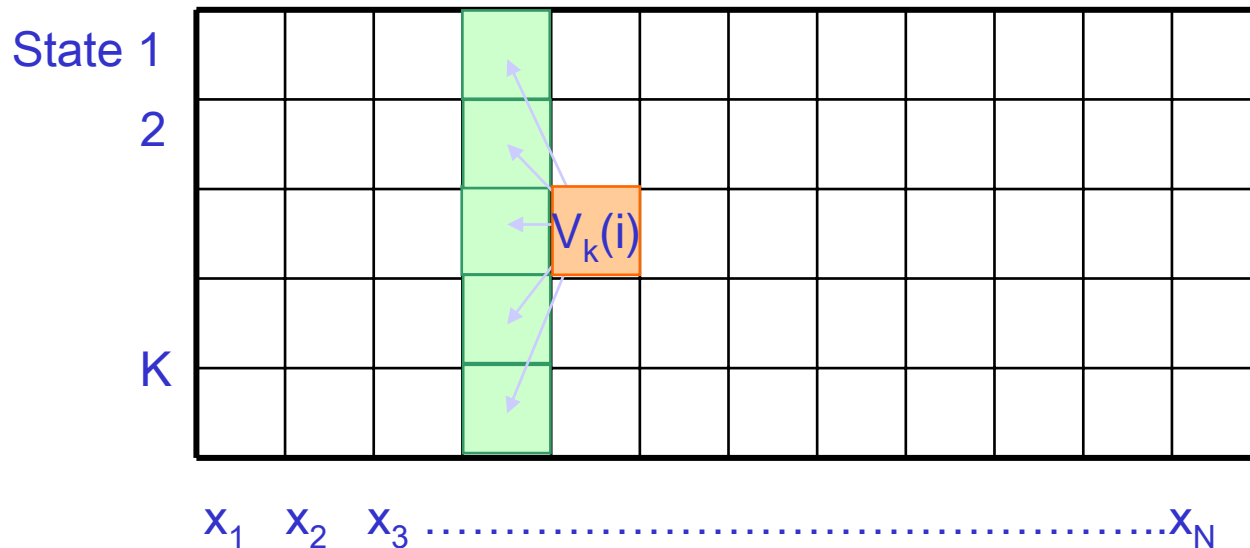
GHMMs as Generative Models

Like HMM, we can use a GHMM to generate a sequence and state labeling

- Choose initial state, q_1 , from π_i
- Choose a duration d_1 , from length distribution $f_{q_1}(d)$
- Generate a subsequence from state model $e_{q_1}(X_{1,d})$
- Choose a subsequent state, q_2 , conditioned on the current state, according to $a_{jk} = P(q_k | q_j)$
- Repeat until number of nucleotides equals desired length of sequence

Given a sequence, how do we pick a state labeling (segmentation)?

The Viterbi Algorithm - HMMs



Input: $x = x_1 \dots x_N$

Initialization:

$$V_0(0)=1, V_k(0) = 0, \text{ for all } k > 0$$

Iteration:

$$V_k(i) = e_k(x_i) \times \max_j a_{jk} V_j(i-1)$$

Termination:

$$P(x, \pi^*) = \max_k V_k(N)$$

Traceback:

Follow max pointers back

In practice:

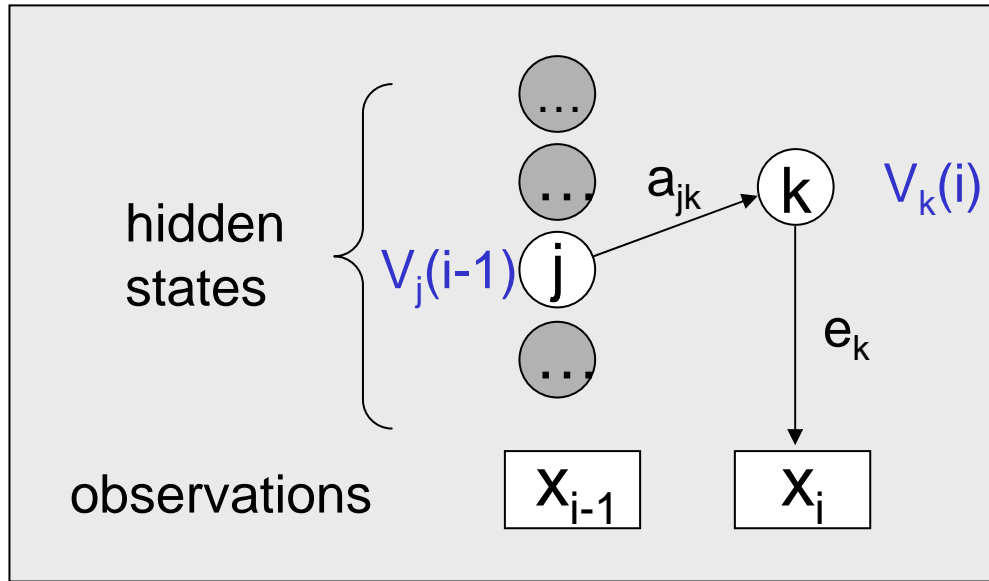
Use log scores for computation

Running time and space:

Time: $O(K^2N)$

Space: $O(KN)$

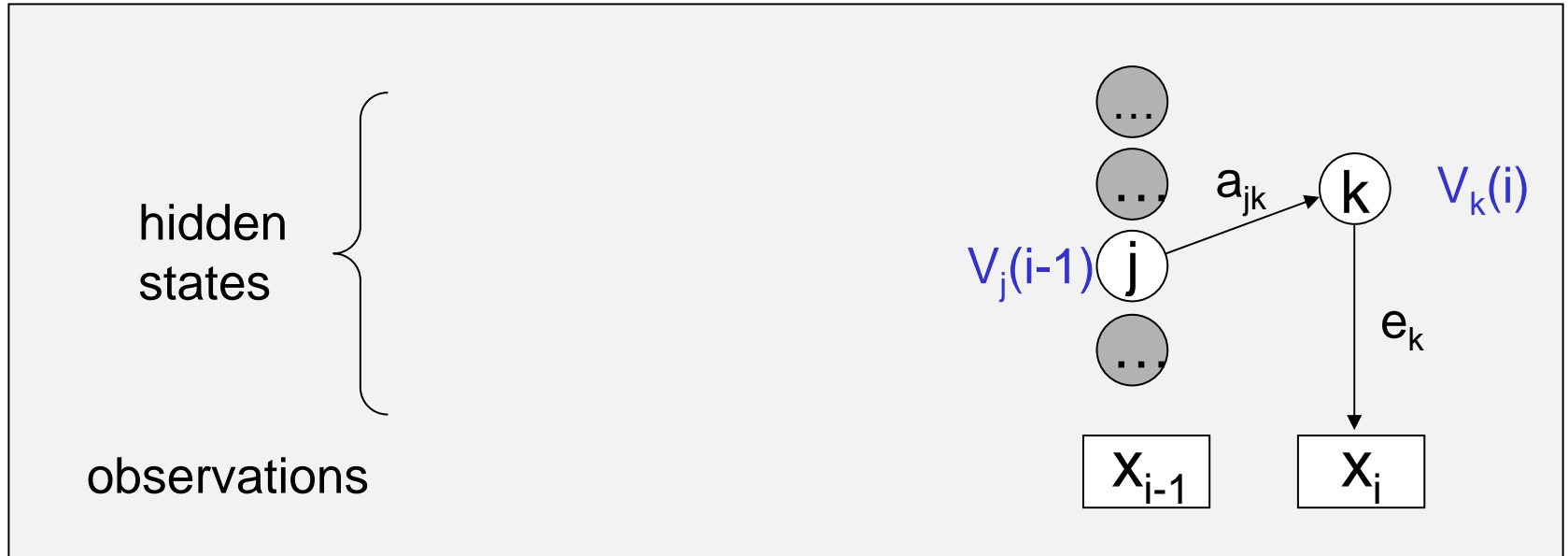
HMM Viterbi Recursion



- Assume we know V_j for the previous time step (i-1)

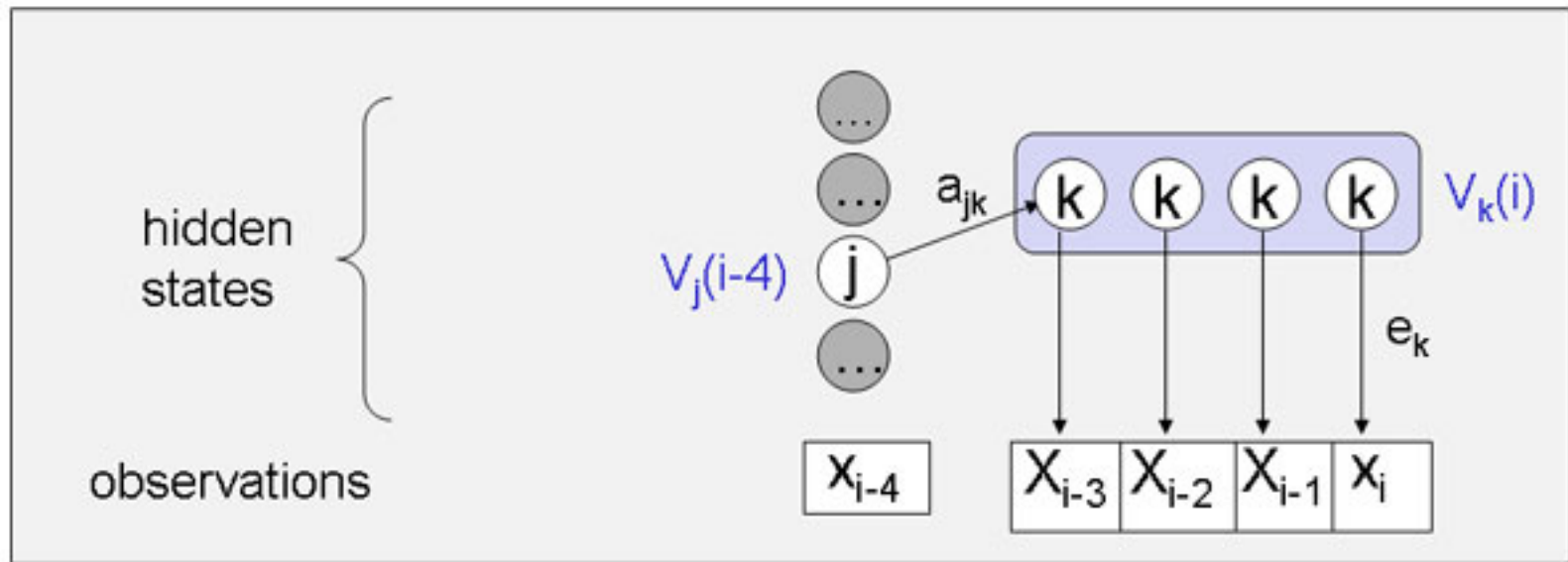
- Calculate $V_k(i) = e_k(x_i) * \max_j (V_j(i-1) \times a_{jk})$
- current max
this emission
max ending in state j at step i-1
Transition from state j
- all possible previous states j

GHMM Recursion



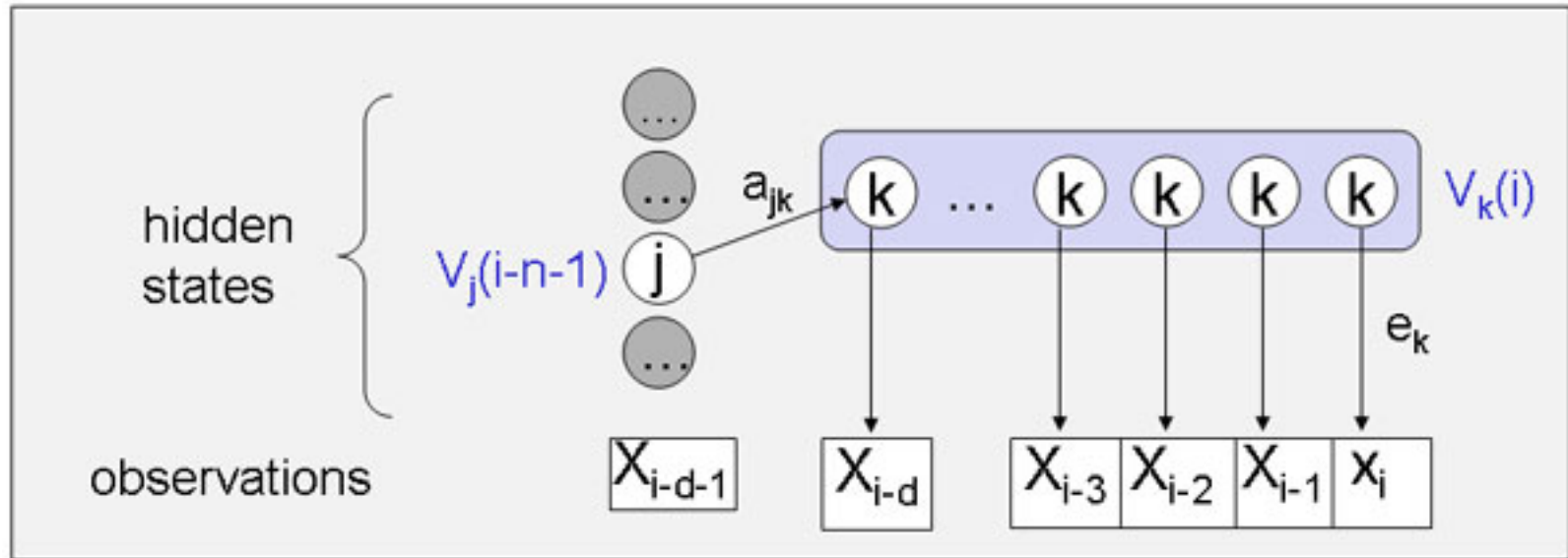
We could have come from state j at the last position...

GHMM Recursion



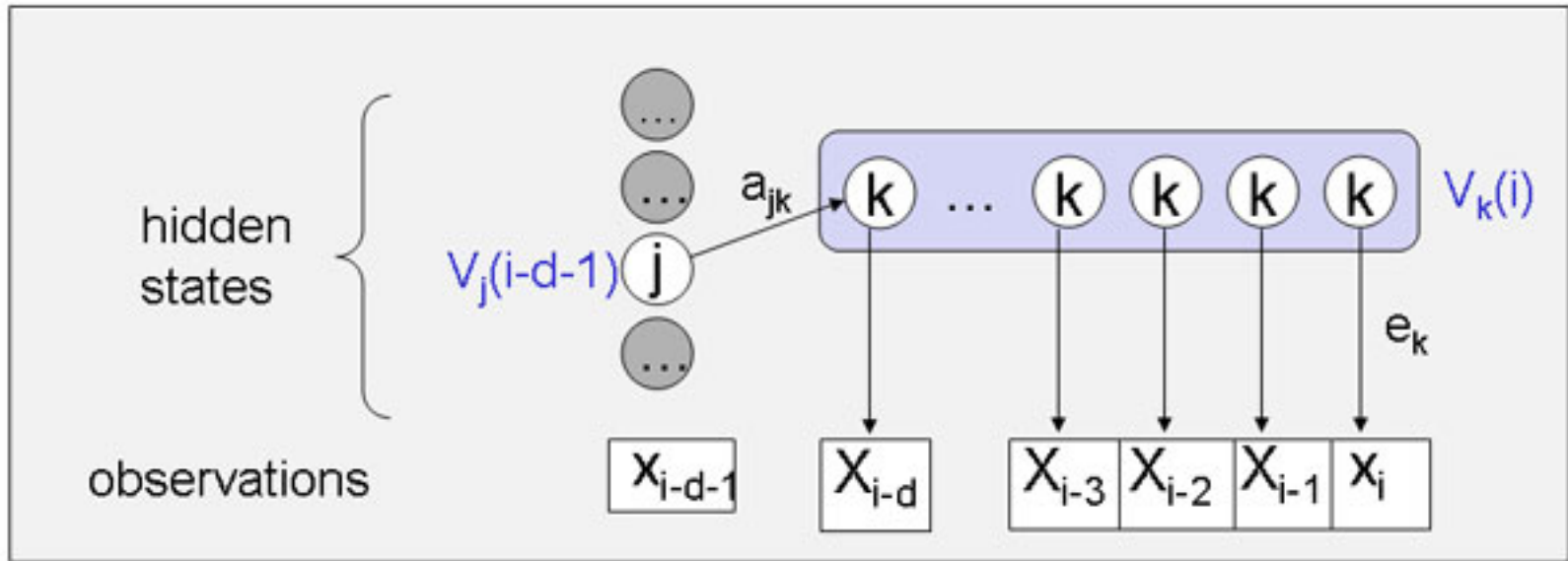
But we could also have come from state j
4 positions ago...
(State k could have duration 4 so far)

GHMM Recursion



In general, we could have come from state j
d positions ago
(State k could have duration d so far)

GHMM Recursion



This leads to the following recursion equation:

$$V_k(i) = \max_j \max_d \left[P(x_i \dots x_{i-d} | k) \cdot V_j(i-d) \cdot P(d|k) \cdot a_{jk} \right]$$

Max over all prev
states *and* state
durations

Prob of
subsequence
given state k

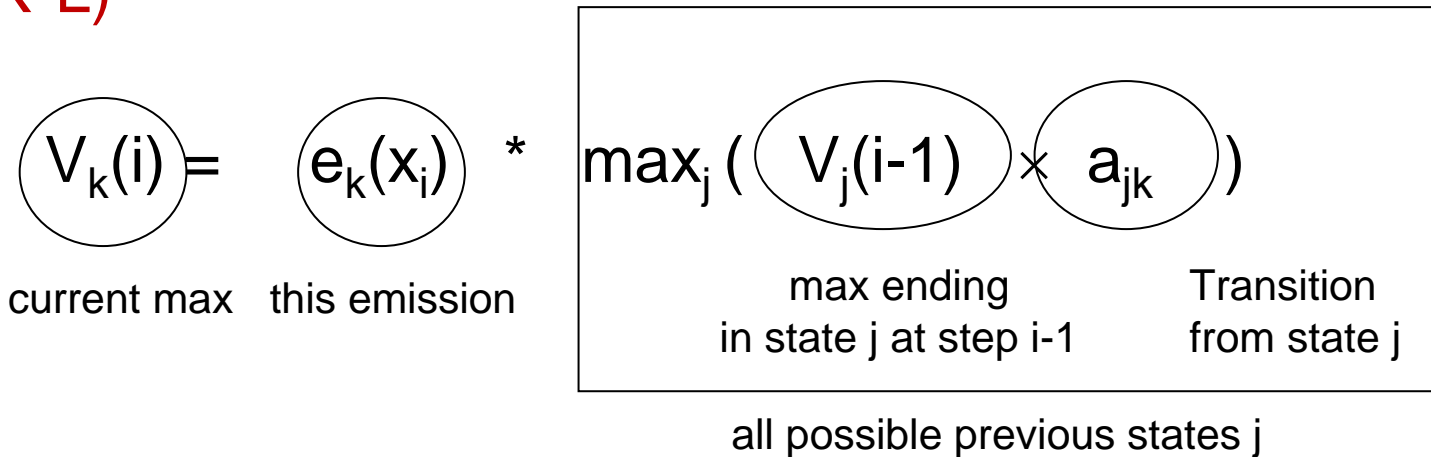
Max ending in
state j at step $i-d$

Prob that state
 k has duration
 d

Transition
from $j \rightarrow k$

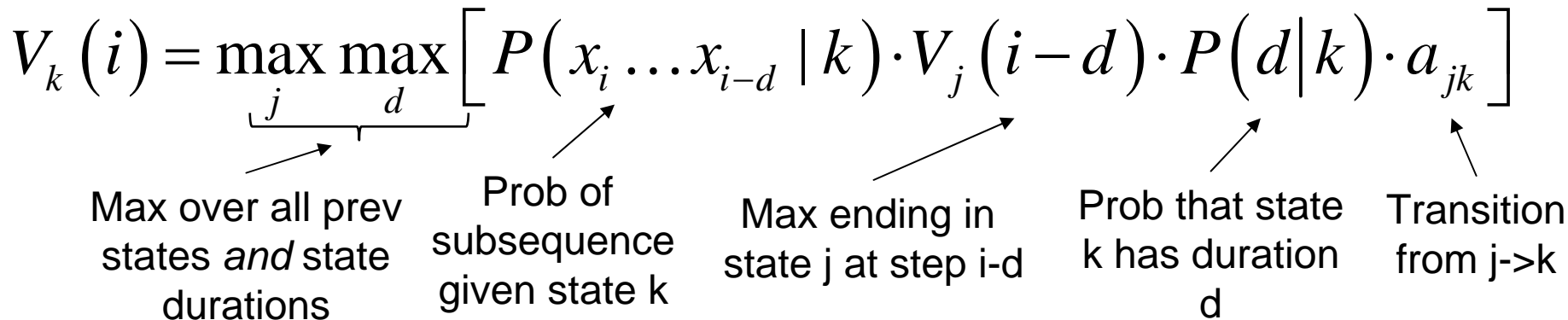
Comparing GHMMs and HMMs

HMM – $O(K^2L)$



GHMM – $O(K^2L^3)$

Similar modifications needed for forward and backward algorithms



Prediction of Complete Gene Structures in Human Genomic DNA

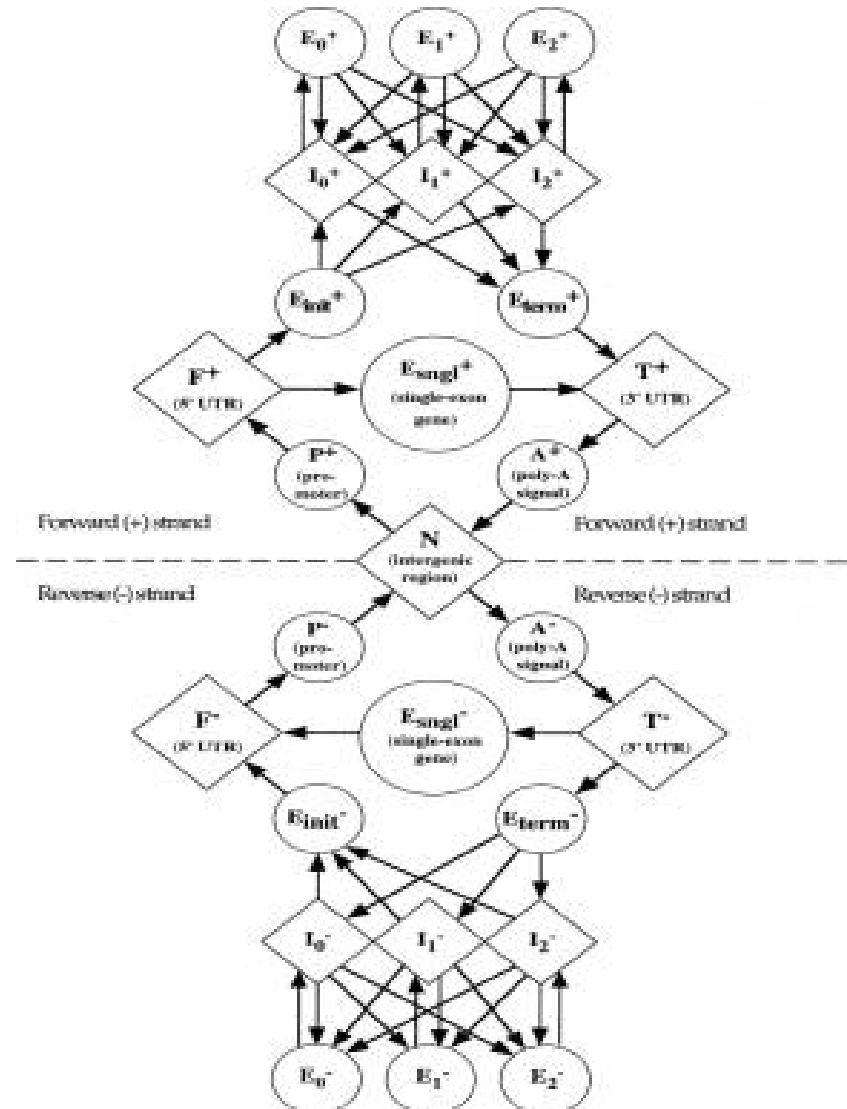
Chris Burge* and Samuel Karlin

*Department of Mathematics
Stanford University, Stanford
CA, 94305, USA*

We introduce a general probabilistic model of the gene structure of human genomic sequences which incorporates descriptions of the basic transcriptional, translational and splicing signals, as well as length distributions and compositional features of exons, introns and intergenic regions. Distinct sets of model parameters are derived to account for the many substantial differences in gene density and structure observed in distinct C + G compositional regions of the human genome. In addition, new models of the donor and acceptor splice signals are described which capture potentially important dependencies between signal positions. The model is applied to the problem of gene identification in a computer program, GENSCAN, which identifies complete exon/intron structures of genes in genomic DNA. Novel features of the program include the capacity to predict multiple genes in a sequence, to deal with partial as well as complete genes, and to predict consistent sets of genes occurring on either or both DNA strands. GENSCAN is shown to have substantially higher accuracy than existing methods when tested on standardized sets of human and vertebrate genes, with 75 to 80% of exons identified exactly. The program is also capable of indicating fairly accurately the reliability of each predicted exon. Consistently high levels of accuracy are observed for sequences of differing C + G content and for distinct groups of vertebrates.

Genscan - Burge and Karlin, (1997)

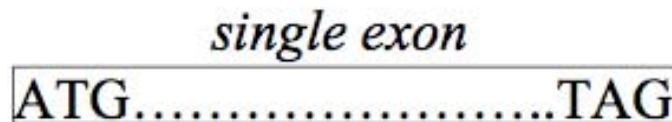
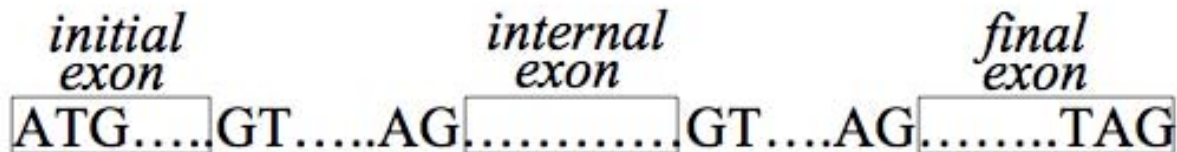
- Explicit State Duration GHMM
- 5th order markov models for coding and non-coding sequences
- Each CDS frame has own model
- WAM models for start/stop codons and acceptor sites
- MDD model for donor sites
- Separate parameters for regions of different GC content
- Model +/- strand concurrently



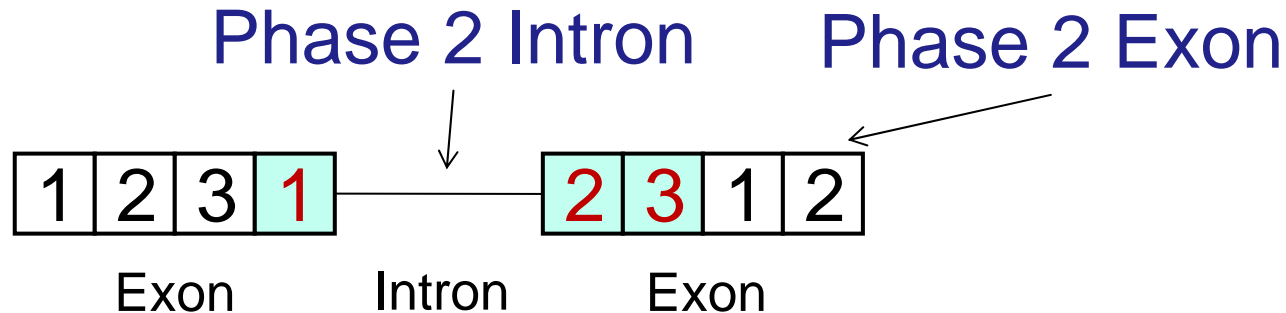
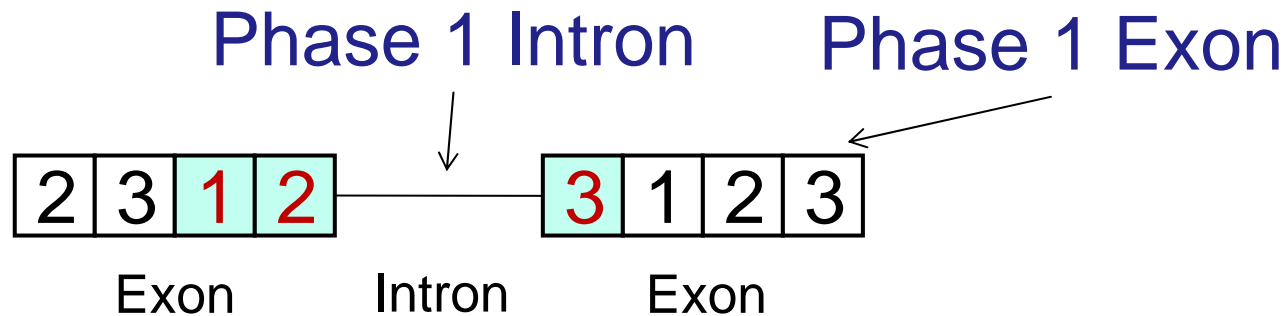
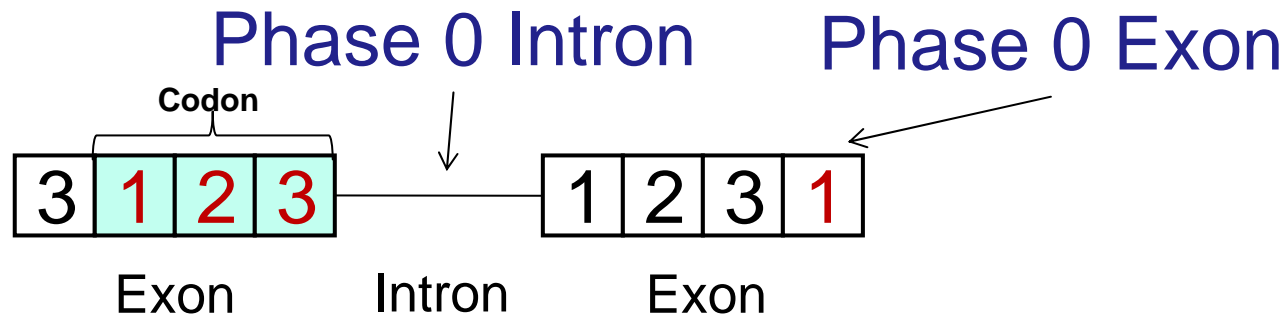
Types of Exons

Three types of exons are defined, for convenience:

- *initial exons* extend from a start codon to the first donor site;
- *internal exons* extend from one acceptor site to the next donor site;
- *final exons* extend from the last acceptor site to the stop codon;
- *single exons* (which occur only in *intronless genes*) extend from the start codon to the stop codon:



Intron and Exon Phase

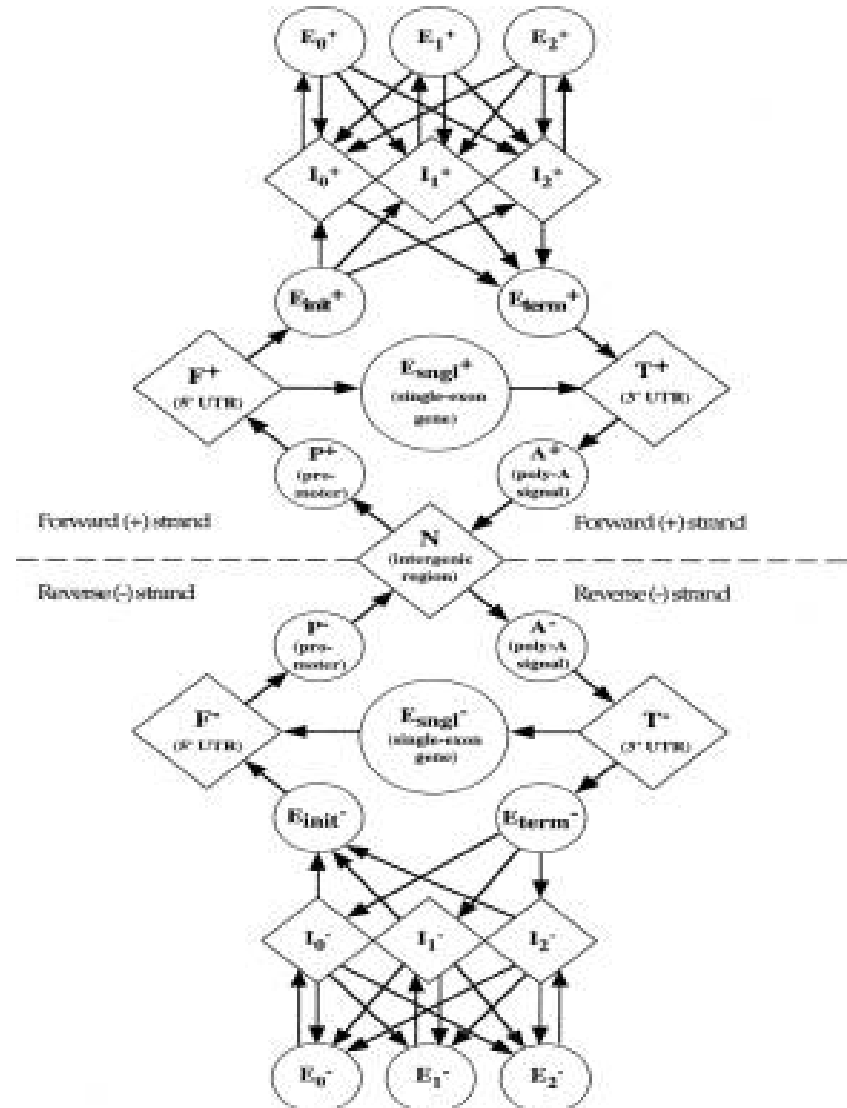


Two State Types in Genscan

- D-type represented by diamonds
- C-type represented by circles

D states are always followed by C and vice versa

C states always preceded by same D-state



Two State Types in Genscan

- **D-Type** – geometric length distribution
 - Intergenic regions
 - UTRs
 - Introns

Sequence models are “factorable”:

$$P_k(Xa, c) = P_k(Xa, b)P_k(Xb + 1, c)$$

$$f_k(d) = P(\text{state duration } d | \text{state } k) = p_k f_k(d - 1)$$

Increasing duration by one changes probability by constant factor

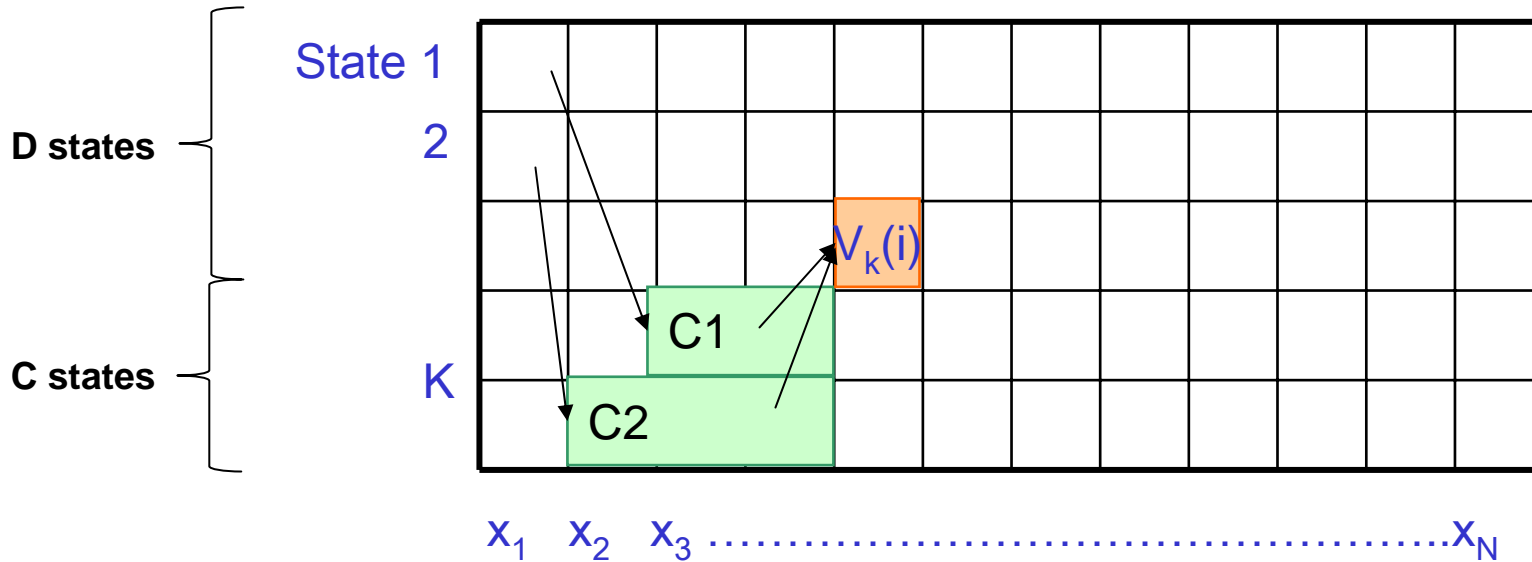
Two State Types in Genscan

- **C-Type** – general length distributions and sequence generating modes
 - Exons (initial, internal, terminal)
 - Promotors
 - Poly-A Signals

Genscan Inference

- Genscan uses same basic forward, backward, and viterbi algorithms as generic GHMMs
- But assumptions about C, and D states reduce algorithmic complexity

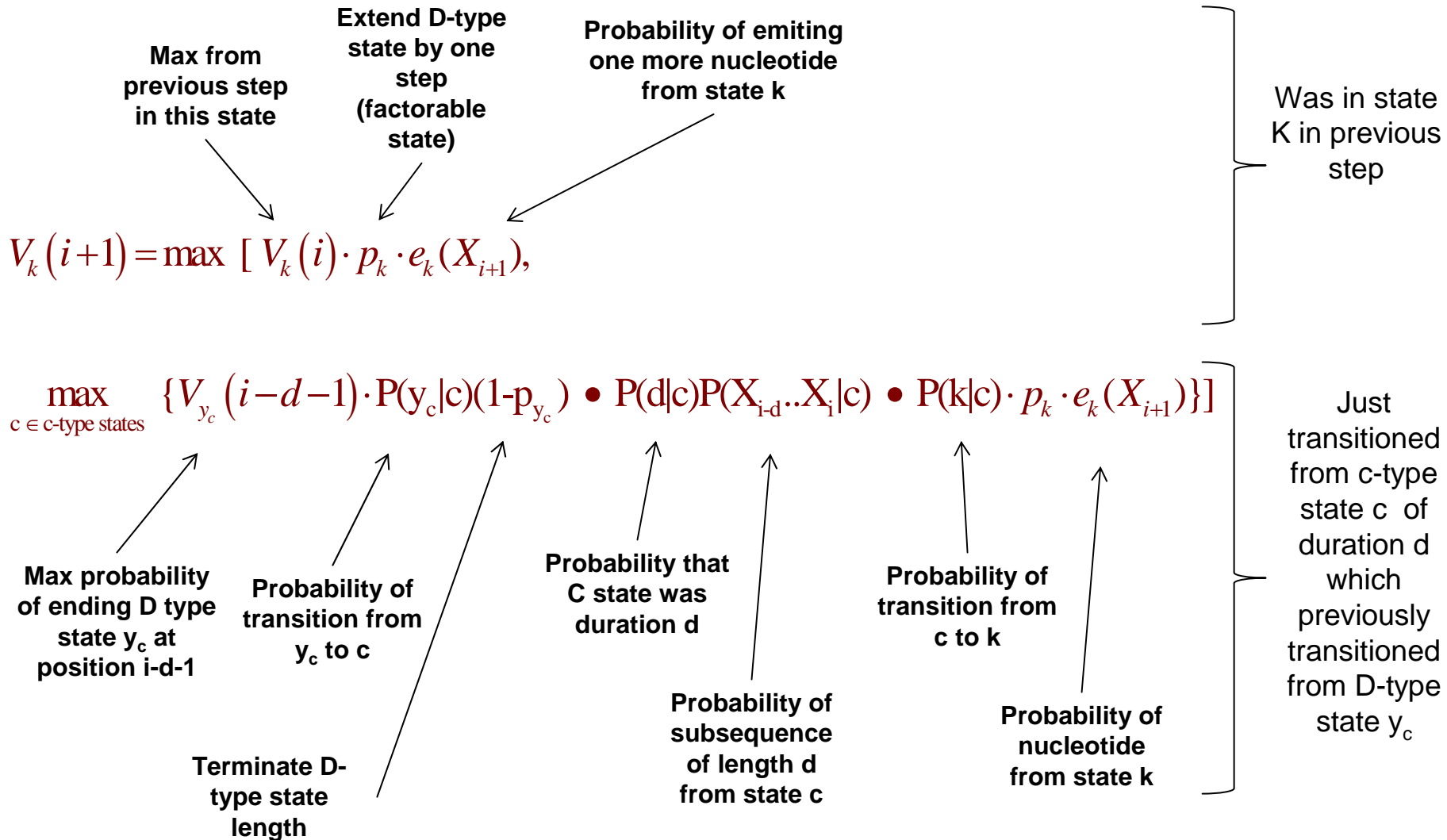
Genscan Inference – C-state List



$L_k(i)$

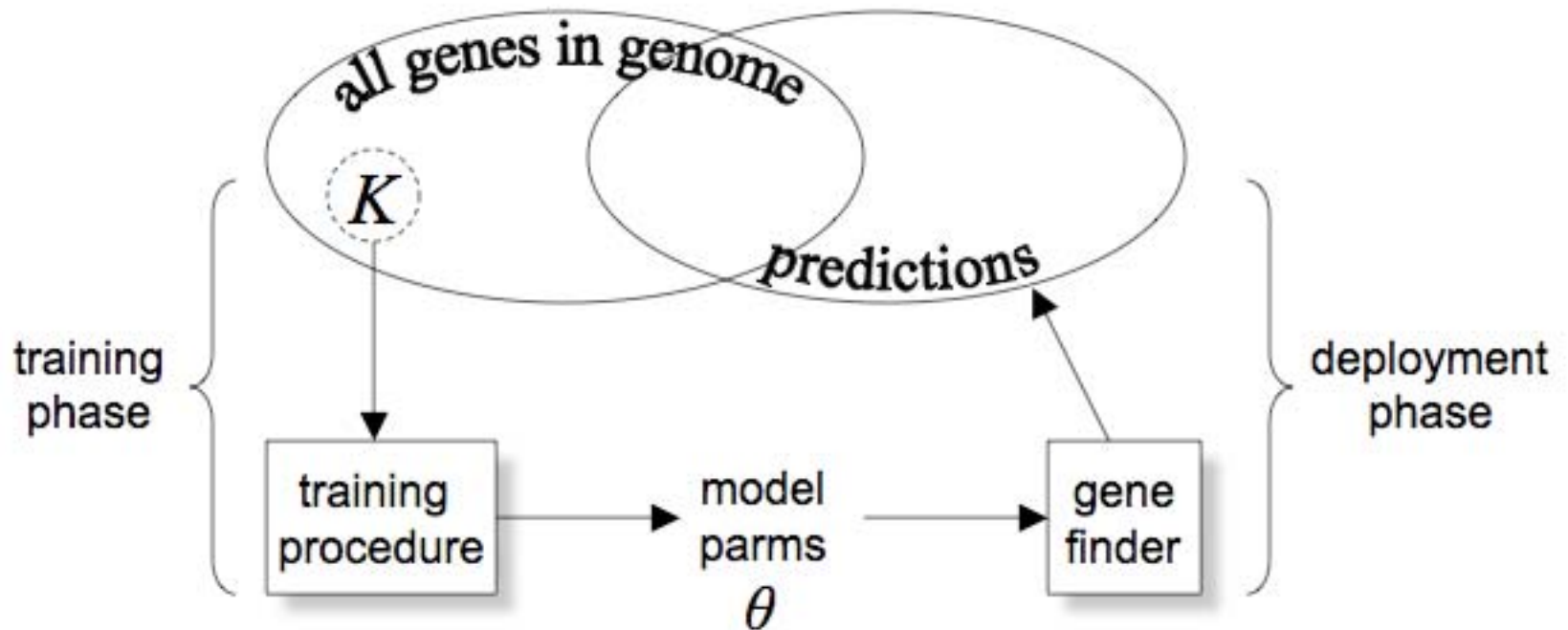
C1, duration 2, previous D-state=1
C2, duration 3, previous D-state=2

Genscan Viterbi Induction



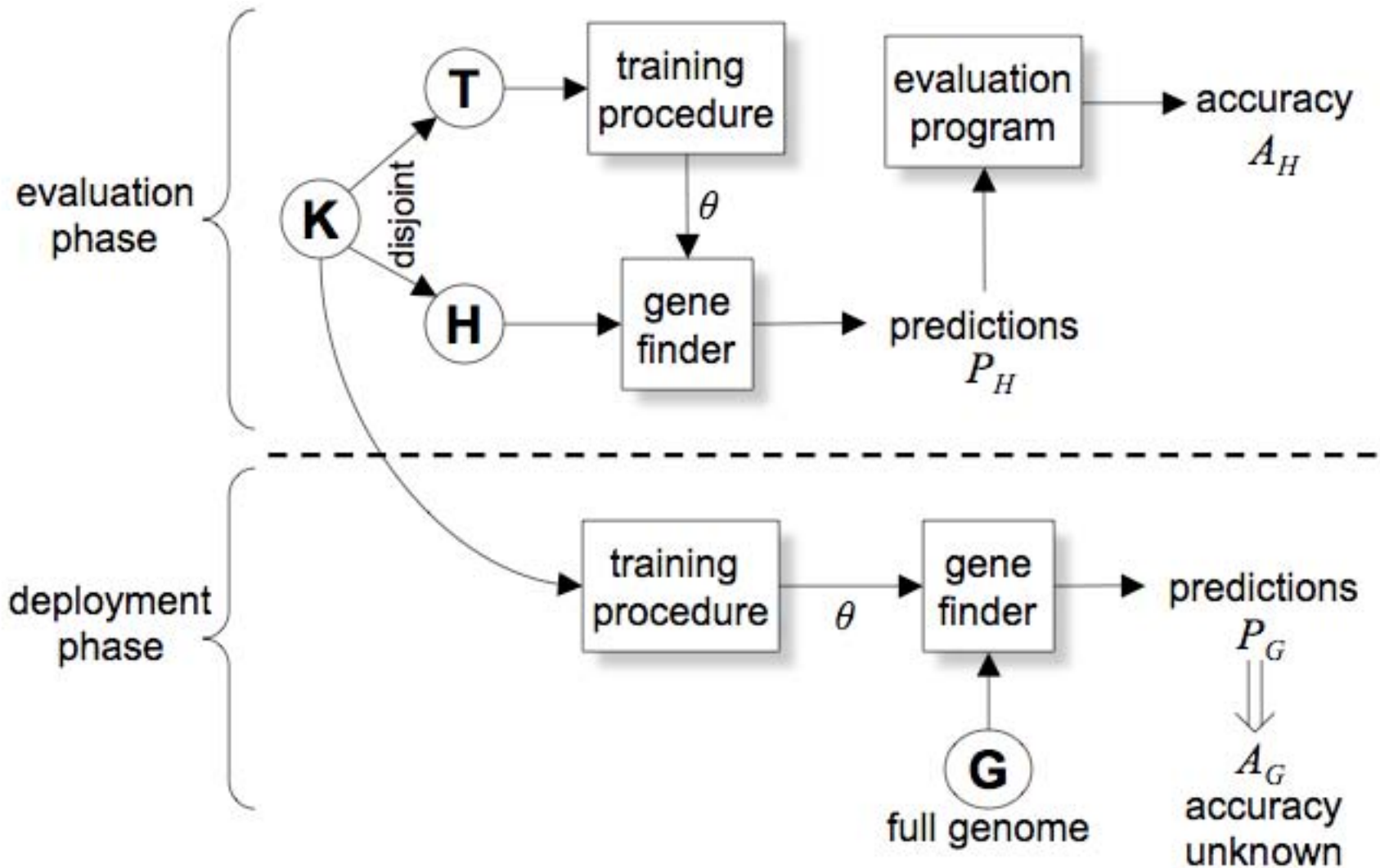
Training A Gene Predictor

During *training* of a gene finder, only a subset K of an organism's gene set will be available for training:



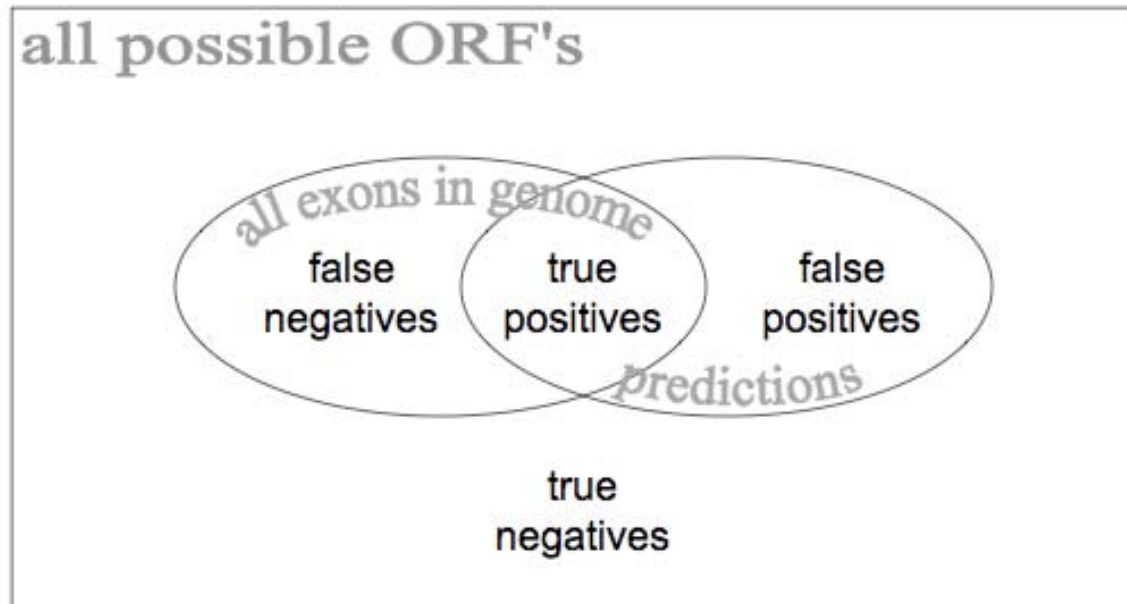
The gene finder will later be *deployed* for use in predicting the rest of the organism's genes. The way in which the *model parameters* are inferred during training can significantly affect the accuracy of the deployed program.

Training A Gene Predictor



Gene Prediction Accuracy

Gene predictions can be evaluated in terms of *true positives* (predicted features that are real), *true negatives* (non-predicted features that are not real), *false positives* (predicted features that are not real), and *false negatives* (real features that were not predicted):



These definitions can be applied at the *whole-gene*, *whole-exon*, or *individual nucleotide* level to arrive at three sets of statistics.

Accuracy Metrics

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

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

$$F = \frac{2 \times Sn \times Sp}{Sn + Sp}$$

$$SMC = \frac{TP + TN}{TP + TN + FP + FN}$$

$$CC = \frac{(TP \times TN) - (FN \times FP)}{\sqrt{(TP + FN) \times (TN + FP) \times (TP + FP) \times (TN + FN)}}$$

$$ACP = \frac{1}{n} \left(\frac{TP}{TP + FN} + \frac{TP}{TP + FP} + \frac{TN}{TN + FP} + \frac{TN}{TN + FN} \right),$$

$$AC = 2(ACP - 0.5).$$

More Information

- <http://genes.mit.edu/burgelab/links.html>
- <http://www.geneprediction.org/book/classroom.html>