# Applications of HMMs in Computational Biology

BMI/CS 576

www.biostat.wisc.edu/bmi576.html

Mark Craven

craven@biostat.wisc.edu
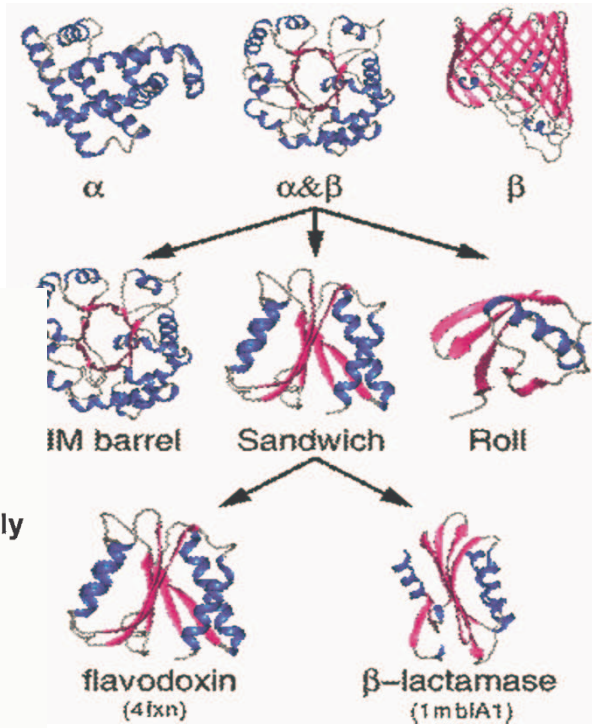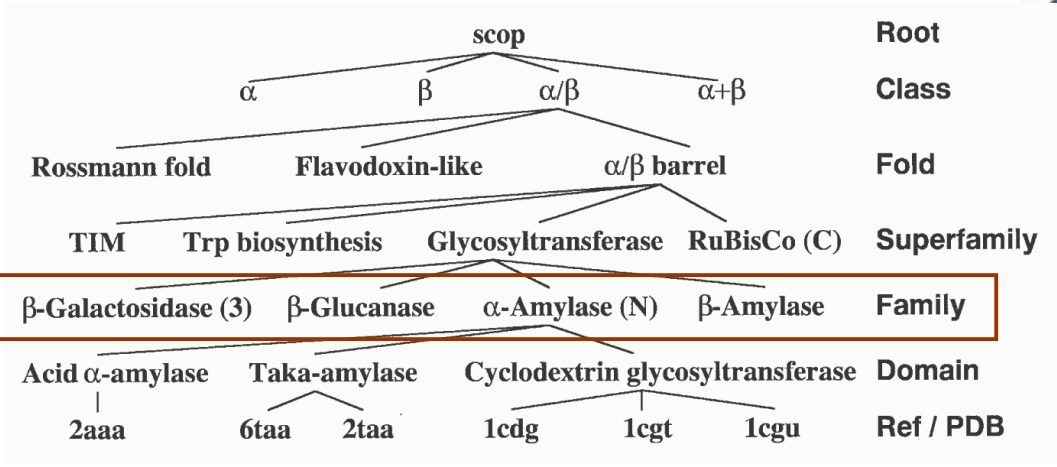
Fall 2011

# The protein classification task

Given: amino-acid sequence of a protein

Do: predict the *family* to which it belongs

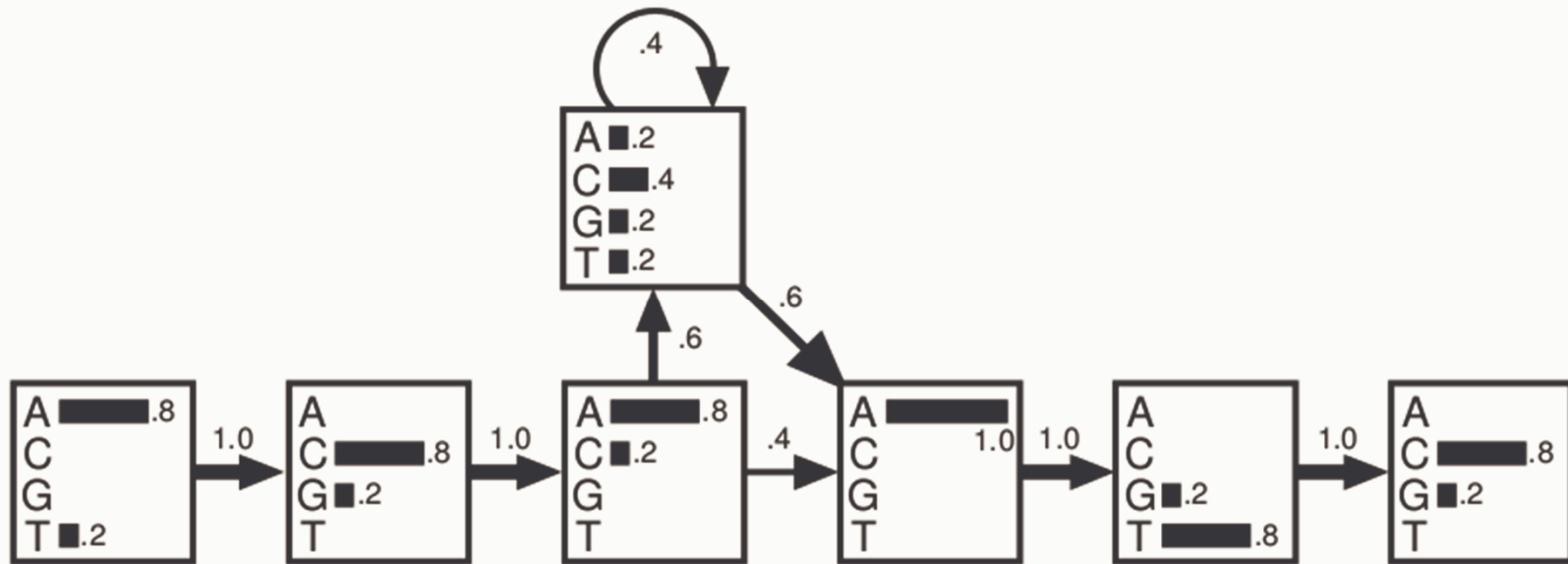GDLSTPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVCVLAHHFGKEFTPPVQAAYAKVVAGVANALAHKYH

# Protein family - a simplified view

```
A C A - - - A T G  ⎤
T C A A C T A T C  ⎥
A C A C - - A G C  ⎬  family
A G A - - - A T C  ⎥
A C C G - - A T C  ⎦
```

```
A C A C - - A T C      query 1
A A A C - - A T C      query 2
T G C T - - A T C      query 3
```

An example from Krogh: An Introduction to HMMs for Biological Sequences, CMMB 1998.

# Protein family - HMM



| | Sequence | Probability ×100 | Log odds |
|---|---|---|---|
| Consensus | A C A C - - A T C | 4.7 | 6.7 |
| Original | A C A - - - A T G | 3.3 | 4.9 |
| sequences | T C A A C T A T C | 0.0075 | 3.0 |
| | A C A C - - A G C | 1.2 | 5.3 |
| | A G A - - - A T C | 3.3 | 4.9 |
| | A C C G - - A T C | 0.59 | 4.6 |
| Exceptional | T G C T - - A G G | 0.0023 | -0.97 |

An example from Krogh: An Introduction to HMMs for Biological Sequences, CMMB 1998.
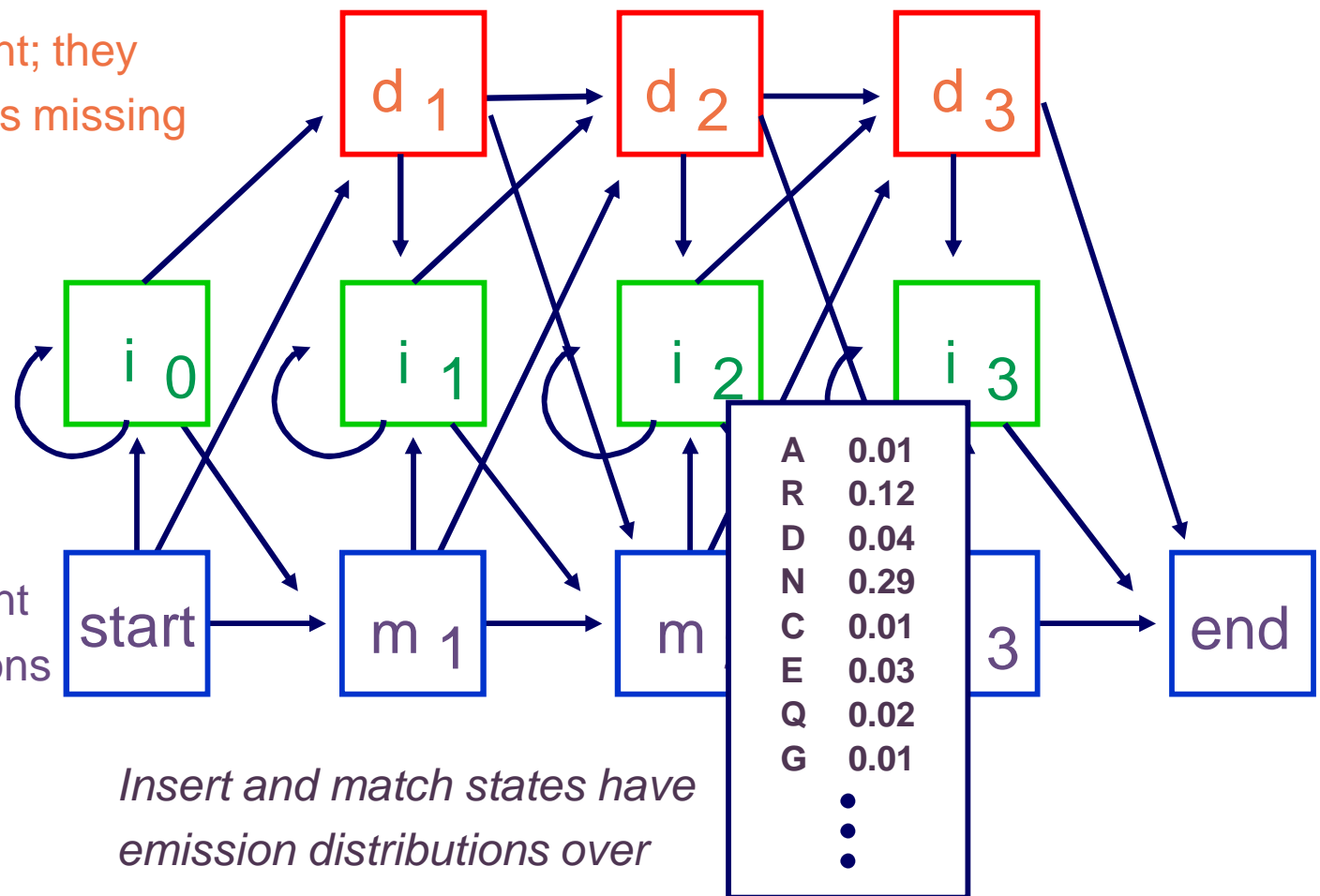
# Profile HMMs

- profile HMMs are used to model families of sequences

*Delete states* are silent; they Account for characters missing in some sequences

*Insert states* account for extra characters in some sequences

*Match states* represent key conserved positions

| | | |
|---|---|---|
| A | 0.01 |
| R | 0.12 |
| D | 0.04 |
| N | 0.29 |
| C | 0.01 |
| E | 0.03 |
| Q | 0.02 |
| G | 0.01 |

$d_1$  $d_2$  $d_3$

$i_0$  $i_1$  $i_2$  $i_3$

start  $m_1$  $m$  $3$  end

*Insert and match states have emission distributions over sequence characters*

# Multiple alignment of SH3 domain

```
G G W W R G d y . g g k k q L W F P S N Y V
I G W L N G y n e t t g e r G D F P G T Y V
P N W W E G q l . . n n r r G I F P S N Y V
D E W W Q A r r . . d e q i G I V P S K - -
G E W W K A q s . . t g q e G F I P F N F V
G D W W L A r s . . s g q t G Y I P S N Y V
G D W W D A e l . . k g r r G K V P S N Y L
- D W W E A r s l s s g h r G Y V P S N Y V
G D W W Y A r s l i t n s e G Y I P S T Y V
G E W W K A r s l a t r k e G Y I P S N Y V
G D W W L A r s l v t g r e G Y V P S N F V
G E W W K A k s l s s k r e G F I P S N Y V
G E W C E A q t . k n g q . G W V P S N Y I
S D W W R V v n l t t r q e G L I P L N F V
L P W W R A r d . k n g q e G Y I P S N Y I
R D W W E F r s k t v y t p G Y Y E S G Y V
E H W W K V k d . a l g n v G Y I P S N Y V
I H W W R V q d . r n g h e G Y V P S S Y L
K D W W K V e v . . n d r q G F V P A A Y V
V G W M P G l n e r t r q r G D F P G T Y V
P D W W E G e l . . n g q r G V F P A S Y V
E N W W N G e i . . g n r k G I F P A T Y V
E E W L E G e c . . k g k v G I F P K V F V
G G W W K G d y . g t r i q Q Y F P S N Y V
D G W W R G s y . . n g q v G W F P S N Y V
Q G W W R G e i . . y g r v G W F P A N Y V
G R W W K A r r . a n g e t G I I P S N Y V
G G W T Q G e l . k s g q k G W A P T N Y L
G D W W E A r s n . t g e n G Y I P S N Y V
N D W W T G r t . . n g k e G I F P A N Y V
```
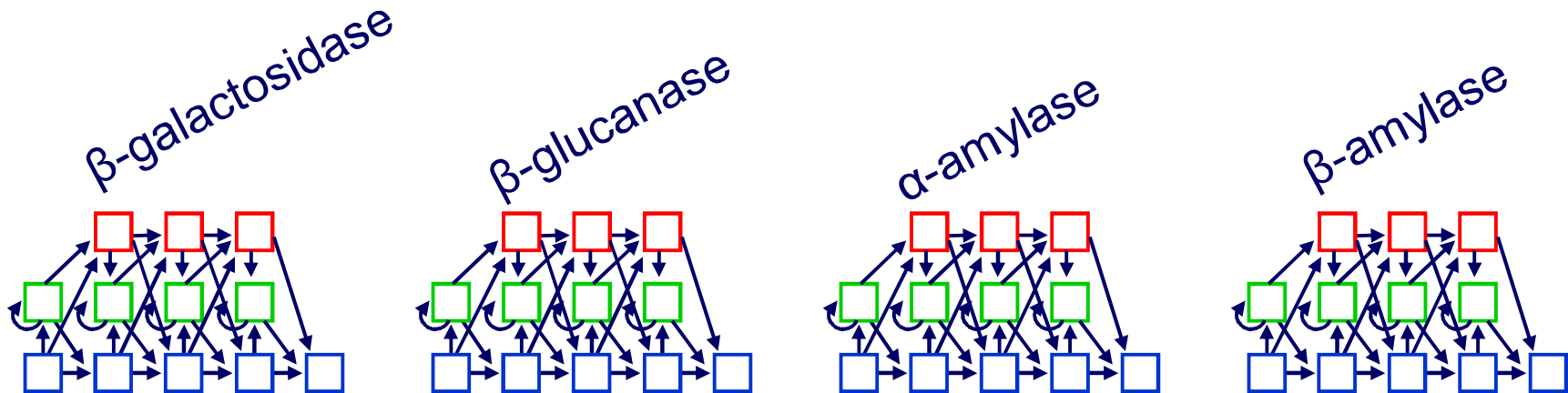
Figure from A. Krogh, An Introduction to Hidden Markov Models for Biological Sequences

# A profile HMM trained for the SH3 domain

insert states

delete states (silent)

match states

# Profile HMMs

- to classify sequences according to family, we can train a profile HMM to model the proteins of each family of interest

- given a sequence $x$, use Bayes' rule to make classification

$$P(c_i \mid x) = \frac{P(x \mid c_i)P(c_i)}{\displaystyle\sum_j P(x \mid c_j)P(c_j)}$$

- use Forward algorithm to compute $P(x \mid c_i)$ for each family $c_i$

β-galactosidase    β-glucanase    α-amylase    β-amylase

# Profile HMM accuracy



Figure from Jaakola et al., ISMB 1999

- classifying 2447proteins into 33 families
- *x*-axis represents the median # of negative sequences that score as high as a positive sequence for a given family's model

# See Pfam database for a large collection profile HMMs

# The gene finding task

Given: an uncharacterized DNA sequence

Do: locate the genes in the sequence, including the coordinates of individual *exons* and *introns*

# Eukaryotic gene structure

# Sources of evidence for gene finding

- **signals**: the sequence *signals* (e.g. splice junctions) involved in gene expression

- **content**: statistical properties that distinguish protein-coding DNA from non-coding DNA

- **conservation**: signal and content properties that are conserved across related sequences (e.g. syntenic regions of the mouse and human genome)

# Gene finding: search by content

- encoding a protein affects the statistical properties of a DNA sequence

| | | | |
|---|---|---|---|
| UUU F 0.46 | UCU S 0.19 | UAU Y 0.44 | UGU C 0.46 |
| UUC F 0.54 | UCC S 0.22 | UAC Y 0.56 | UGC C 0.54 |
| UUA L 0.08 | UCA S 0.15 | UAA * 0.30 | UGA * 0.47 |
| UUG L 0.13 | UCG S 0.05 | UAG * 0.24 | UGG W 1.00 |
| | | | |
| CUU L 0.13 | CCU P 0.29 | CAU H 0.42 | CGU R 0.08 |
| CUC L 0.20 | CCC P 0.32 | CAC H 0.58 | CGC R 0.18 |
| CUA L 0.07 | CCA P 0.28 | CAA Q 0.27 | CGA R 0.11 |
| CUG L 0.40 | CCG P 0.11 | CAG Q 0.73 | CGG R 0.20 |
| | | | |
| AUU I 0.36 | ACU T 0.25 | AAU N 0.47 | AGU S 0.15 |
| AUC I 0.47 | ACC T 0.36 | AAC N 0.53 | AGC S 0.24 |
| AUA I 0.17 | ACA T 0.28 | AAA K 0.43 | AGA R 0.21 |
| AUG M 1.00 | ACG T 0.11 | AAG K 0.57 | AGG R 0.21 |
| | | | |
| GUU V 0.18 | GCU A 0.27 | GAU D 0.46 | GGU G 0.16 |
| GUC V 0.24 | GCC A 0.40 | GAC D 0.54 | GGC G 0.34 |
| GUA V 0.12 | GCA A 0.23 | GAA E 0.42 | GGA G 0.25 |
| GUG V 0.46 | GCG A 0.11 | GAG E 0.58 | GGG G 0.25 |

[Codon/a.a./fraction per codon per a.a.]
Homo sapiens data from the Codon Usage Database

# The GENSCAN HMM for Eukaryotic Gene Finding [Burge & Karlin '97]

Each shape denotes a functional unit of a gene or genomic region and is represented by a submodel in the HMM

Pairs of intron/exon units represent the different ways an intron can interrupt a coding sequence (after 1st base in codon, after 2nd base or after 3rd base)

Complementary submodel (not shown) detects genes on opposite DNA strand
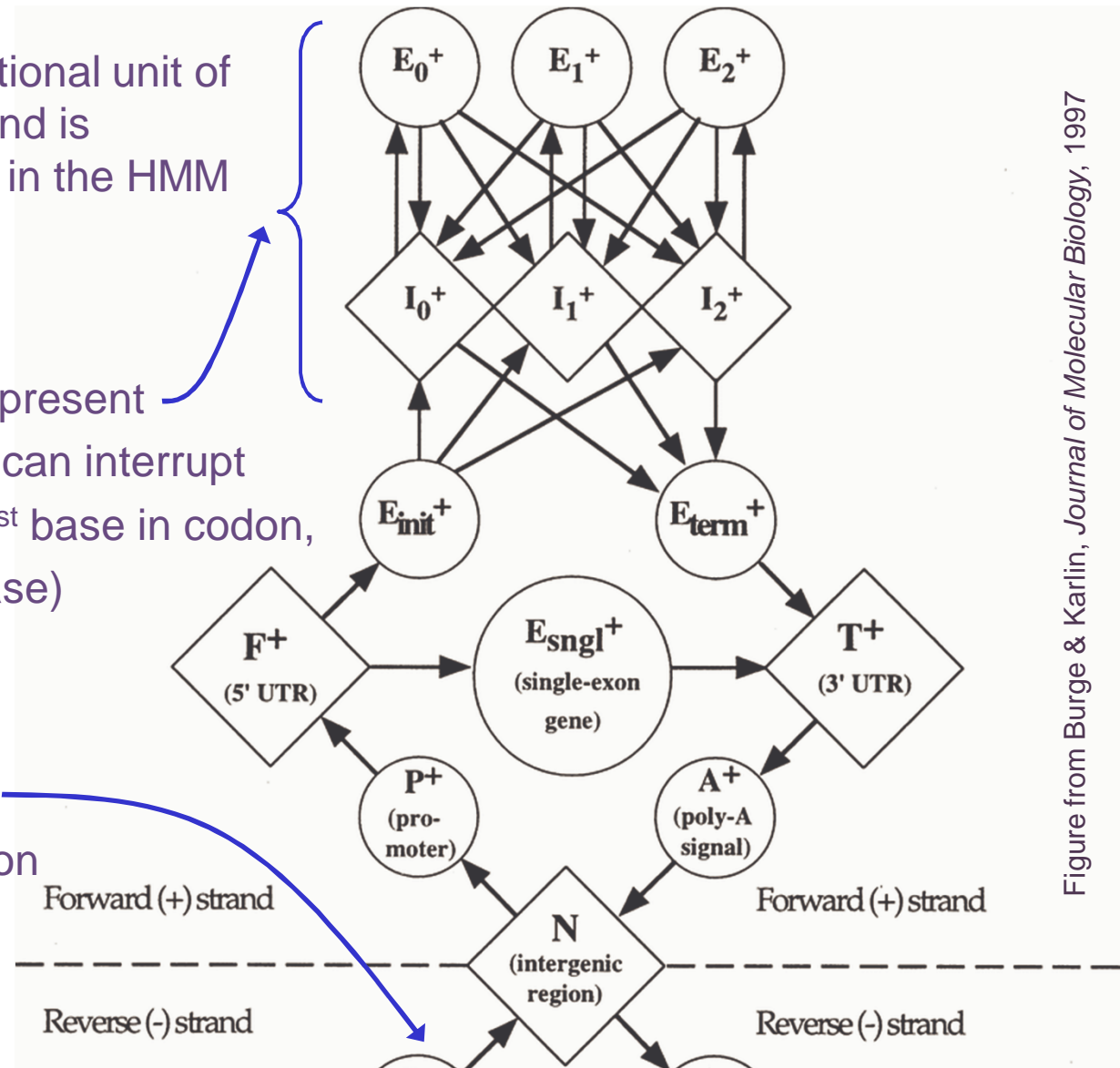


Figure from Burge & Karlin, Journal of Molecular Biology, 1997

# GENSCAN uses a variety of submodel types

| sequence feature | model |
| --- | --- |
| exons | 5$^{th}$ order inhomogenous |
| introns, intergenic regions | 5$^{th}$ order homogenous |
| poly-A, translation initiation, promoter | 0$^{th}$ order, fixed-length |
| splice junctions | tree-structured variable memory |

# Markov models & exons

- consider modeling a given coding sequence
- for each "word" we evaluate, we'll want to consider its position with respect to the reading frame we're assuming

reading frame

G C T A C G G A G C T T C G G A G C

G C T A C **G**    G is in 3rd codon position
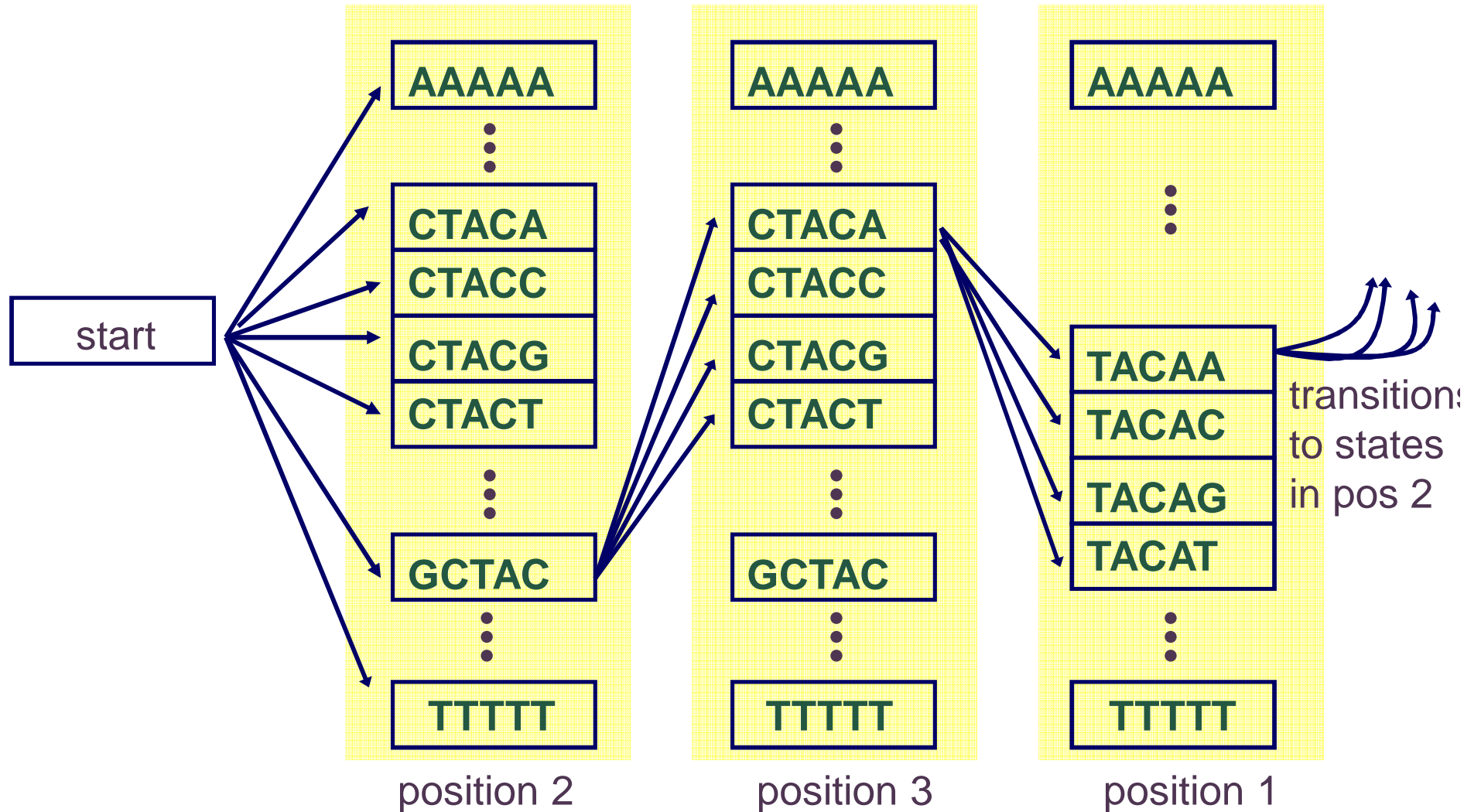
C T A C G **G**    G is in 1st position

T A C G G **A**    A is in 2nd position

- can do this using an inhomogeneous model

# A fifth-order inhomogenous Markov chain



start

| AAAAA | AAAAA | AAAAA |

| CTACA | CTACA | |
| CTACC | CTACC | |
| CTACG | CTACG | TACAA |
| CTACT | CTACT | TACAC |
| | | TACAG |
| GCTAC | GCTAC | TACAT |
| TTTTT | TTTTT | TTTTT |

transitions to states in pos 2

position 2 position 3 position 1
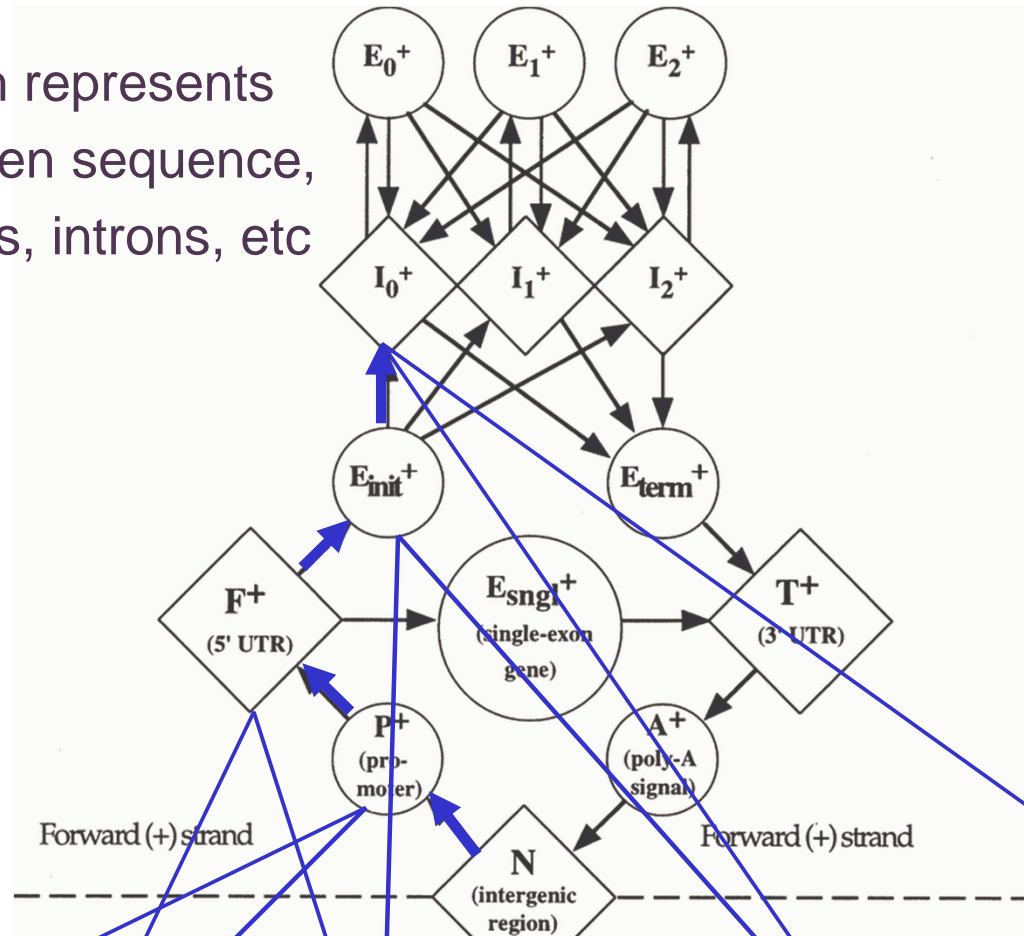
# Inference with the gene-finding HMM

given: an uncharacterized DNA sequence

find: the most probable path through the model for the sequence

- this path will specify the coordinates of the predicted genes (including intron and exon boundaries)
- the Viterbi algorithm is used to compute this path

# Parsing a DNA sequence

The Viterbi path represents
a parse of a given sequence,
predicting exons, introns, etc



ACCGTTACGTGTCATTCTACGTGATCATCGGATCCTAGAATCATCGATCCGTGCGATCGATCGGATTAGCTAGCTTAGCTAGGAGAGCATCGATCGGATCGAGGAGGAGCCTATATAAATCAA

# Other issues in Markov models

- there are many interesting variants and extensions of the models/algorithms we considered here (some of these are covered in BMI/CS 776)
  - separating length/composition distributions with *semi-Markov models*
  - modeling multiple sequences with *pair HMMs*
  - learning the *structure* of HMMs
  - going up the Chomsky hierarchy: *stochastic context free grammars*
  - discriminative learning algorithms (e.g. as in *conditional random fields*)
  - etc.