

# Lecture 17

- Parsing genomes with HMMs

# Genome biology overview

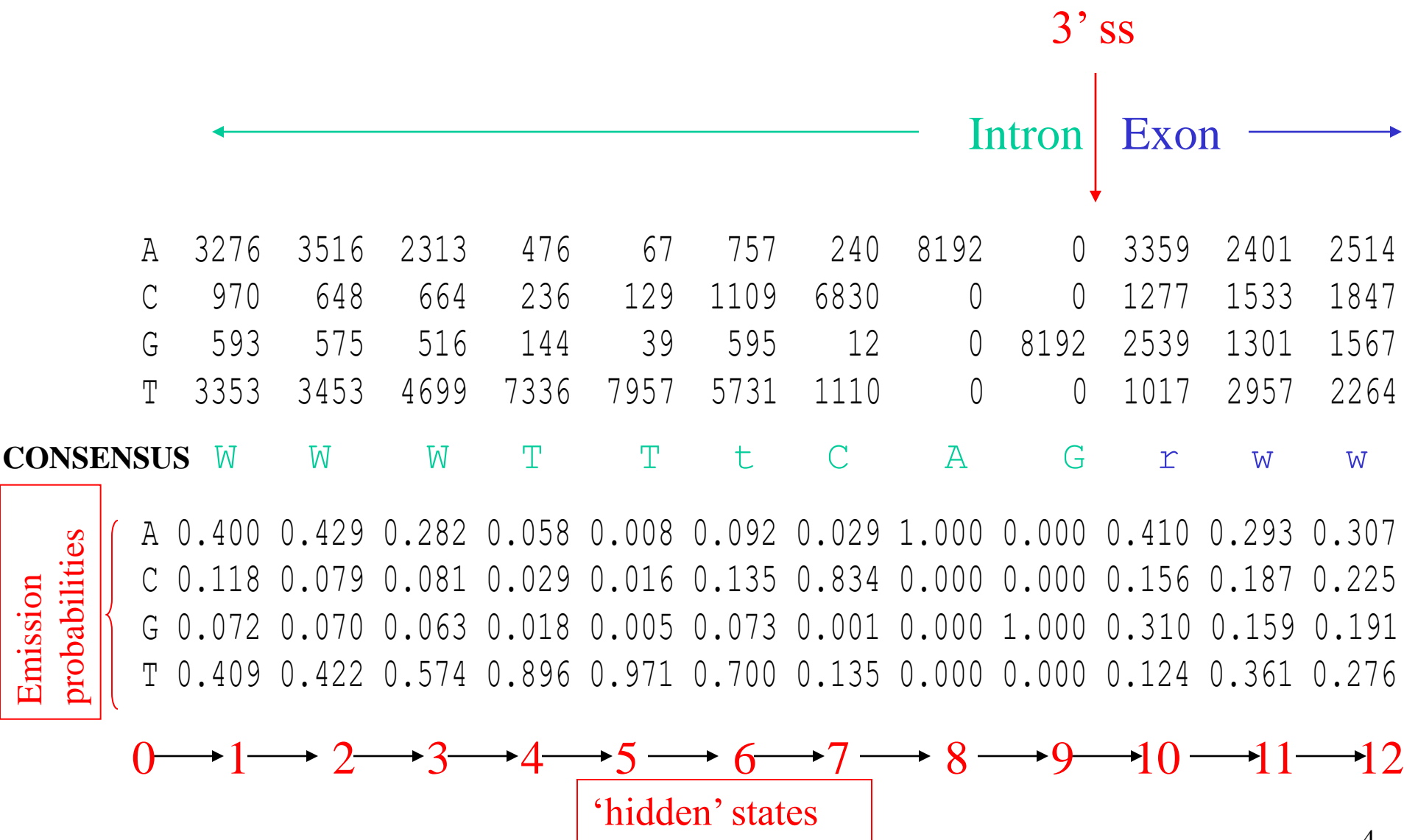
- Genomes undergo two fundamental processes (both involve copying!):
  - Replication
  - Transcription
- Genomic functional information is in the form of *sites*:
  - Short (~2 – ~15 base) sequence segments that bind to an *RNA* or *protein* molecule (the *reader*) to help mediate some function

# Genome HMMs

- a genome consists of (functionally important) *sites* within (nonfunctional) *background* sequence.
- can define an HMM that reflects this:
  - one state per site position, for each type of site
  - background state
  - appropriate topology (allowed transitions)
  - emission & transition probs

and use it to get Viterbi parse & posterior probs

# HMM for *C. elegans* 3' Splice Sites



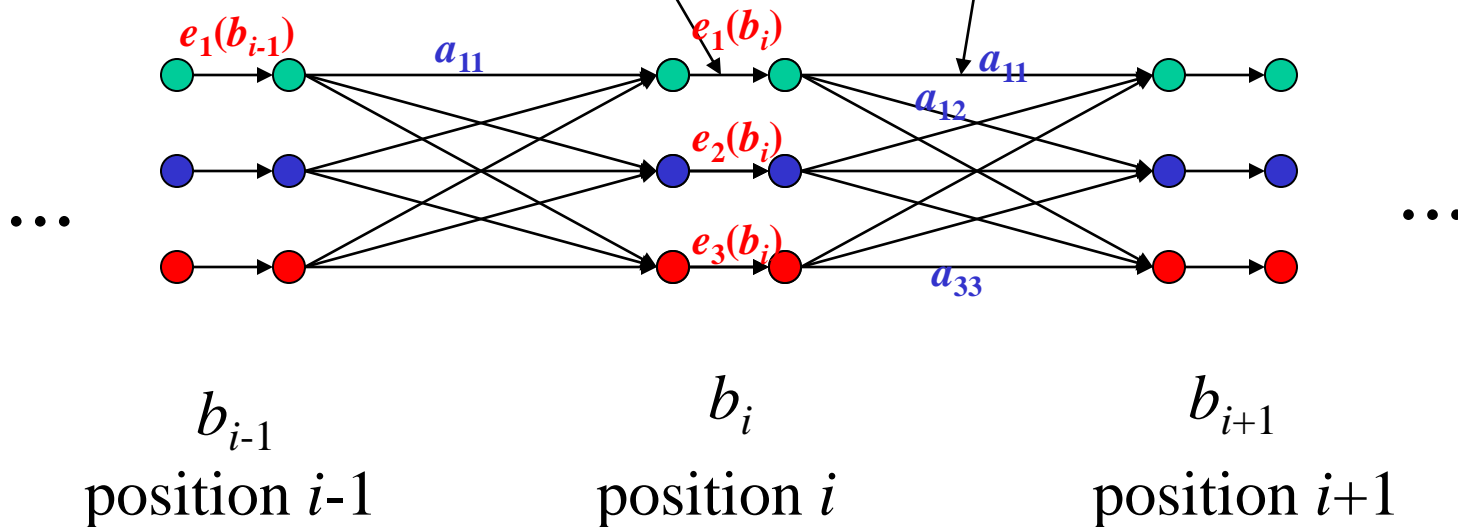
- Complication: sites have *orientation* (top or bottom *strand*)
  - e.g. from transcription direction
- One strategy: analyze 2 strands separately
  - problem: resolving conflicts
- Better strategy: *expand model* to allow sites in both orientations, and run on *top strand* only
  - double # site states
  - bottom strand states have
    - complementary emission probs
    - reversed allowed transitions

- # params does not change
- size of WDAG increases, but only by factor of  $\sim 2$ 
  - no transitions *between* top & bottom strand states, except for background

# *cf. WDAG for 3-state HMM length $n$ sequence (lecture 13)*

weights are emission  
probabilities  $e_k(b_i)$  for  $i^{\text{th}}$   
residue  $b_i$

weights are transition  
probabilities  $a_{kl}$



# Prokaryotes vs eukaryotes

- Such HMMs are most reliable, & most widely used, for prokaryotic genomes, which usually have
  - high site density, homogeneous background
  - relatively simple spatial relationships among sites
  - often relatively little ‘supporting information’ such as
    - protein binding & transcript data
    - closely related genomes



- eukaryotic genomes are less suitable:
  - low site density, heterogeneous background
  - complex site spatial relationships (not well captured by Markov transition model)
  - often much supporting info
    - similar genomes to transfer annotations from;
    - protein binding / RNASeq & other experimental data
    - in principle, some of this could be incorporated into HMM
      - (expanded symbol alphabet)

# Prokaryote genomes

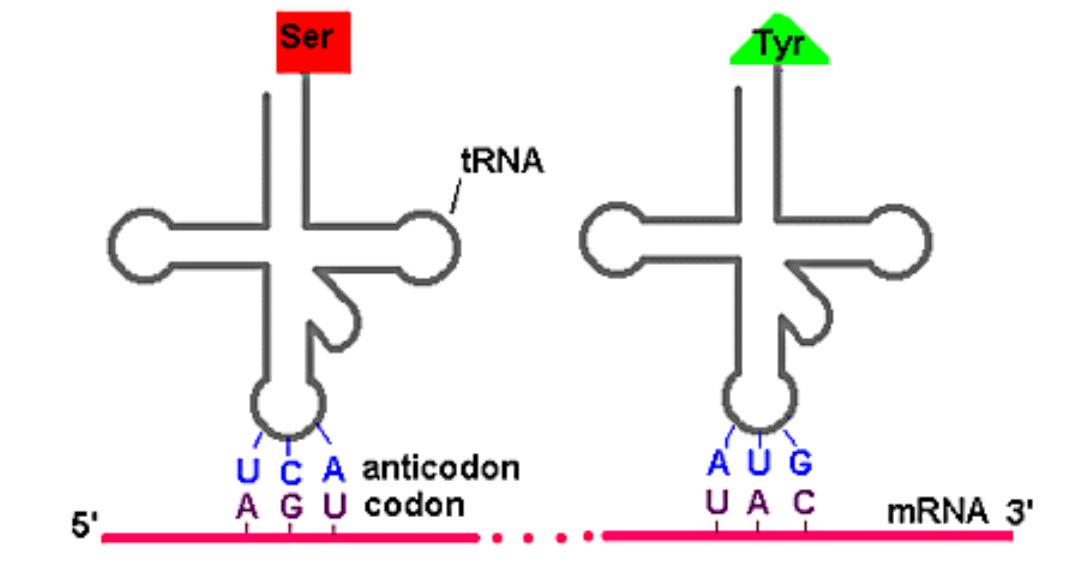
- typically a few MB in size
- up to ~80% protein coding
- Typical CDS size ~1 KB
- introns & overlapping CDSs rare
- range of GC contents

# ORF analysis

- Translate genome in all 6 reading frames
- In each, find ‘open reading frames’ starting with ATG (or NTG), ending in stop
- Sort ORFs by (decreasing) length
- Work through sorted list, discarding any ORF that
  - overlaps a longer one, or
  - is ‘too short’

- Problems:
  - short CDSs are missed
  - CDSs often have long overlapping fake ORFs on opposite strand
  - poor performance on GC-rich genomes (many long fake ORFs)

- Additional information that is present in real coding sequence (but ignored in ORF analysis) – *cf. lecture 3*
  - amino acid usage
  - synonymous codon bias
- Use this, in a probability model!



2nd base in codon

|                   |   | U                        | C                        | A                          | G                         |                  |                   |
|-------------------|---|--------------------------|--------------------------|----------------------------|---------------------------|------------------|-------------------|
| 1st base in codon | U | Phe<br>Phe<br>Leu<br>Leu | Ser<br>Ser<br>Ser<br>Ser | Tyr<br>Tyr<br>STOP<br>STOP | Cys<br>Cys<br>STOP<br>Trp | U<br>C<br>A<br>G | 3rd base in codon |
|                   | C | Leu<br>Leu<br>Leu<br>Leu | Pro<br>Pro<br>Pro<br>Pro | His<br>His<br>Gln<br>Gln   | Arg<br>Arg<br>Arg<br>Arg  | U<br>C<br>A<br>G |                   |
|                   | A | Ile<br>Ile<br>Ile<br>Met | Thr<br>Thr<br>Thr<br>Thr | Asn<br>Asn<br>Lys<br>Lys   | Ser<br>Ser<br>Arg<br>Arg  | U<br>C<br>A<br>G |                   |
|                   | G | Val<br>Val<br>Val<br>Val | Ala<br>Ala<br>Ala<br>Ala | Asp<br>Asp<br>Glu<br>Glu   | Gly<br>Gly<br>Gly<br>Gly  | U<br>C<br>A<br>G |                   |

## The Genetic Code

| Amino Acid | Obs/Exp | 1 <sup>st</sup> codon base | 2 <sup>nd</sup> codon base | 3 <sup>rd</sup> codon base | # codons |
|------------|---------|----------------------------|----------------------------|----------------------------|----------|
| E          | 1.92    | G                          | A                          | R                          | 2        |
| K          | 1.80    | A                          | A                          | R                          | 2        |
| D          | 1.62    | G                          | A                          | Y                          | 2        |
| M          | 1.46    | A                          | T                          | G                          | 1        |
| N          | 1.37    | A                          | A                          | Y                          | 2        |
| F          | 1.25    | T                          | T                          | Y                          | 2        |
| Q          | 1.22    | C                          | A                          | R                          | 2        |
| I          | 1.16    | A                          | T                          | Not G                      | 3        |
| A          | 1.14    | G                          | C                          | N                          | 4        |
| G          | 1.05    | G                          | G                          | N                          | 4        |
| V          | .99     | G                          | T                          | N                          | 4        |
| Y          | .98     | T                          | A                          | Y                          | 2        |
| L          | .95     | C(T)                       | T                          | N                          | 6        |
| T          | .88     | A                          | C                          | N                          | 4        |
| W          | .79     | T                          | G                          | G                          | 1        |
| P          | .74     | C                          | C                          | N                          | 4        |
| S          | .73     | T(A)                       | C(G)                       | N                          | 6        |
| H          | .67     | C                          | A                          | Y                          | 2        |
| R          | .53     | C(A)                       | G                          | N                          | 6        |
| C          | .52     | T                          | G                          | Y                          | 2        |

# Synonymous codon bias

- In most organisms, the codons for an amino acid are not used with equal frequency
- For many organisms this may reflect differences in translational efficiency & accuracy
  - more highly expressed genes have stronger biases
- For mammals codon usage mainly reflects the GC content of the region in which the gene is found
  - GC content variation probably reflects GC-biased gene conversion



|     |         |         |      |         |          |      |         |         |      |         |         |
|-----|---------|---------|------|---------|----------|------|---------|---------|------|---------|---------|
| Phe | 171 UUU | } AAA 0 | Ser  | 147 UCU | } AGA 10 | Tyr  | 124 UAU | } AUA 1 | Cys  | 99 UGU  | } ACA 0 |
|     | 203 UUC |         |      | GAA 14  |          |      | 172 UCC |         |      | GGA 0   |         |
| Leu | 73 UUA  | — UAA 8 | stop | 118 UCA | — UGA 5  | stop | 0 UAA   | — UUA 0 | stop | 0 UGA   | — UCA 0 |
|     | 125 UUG | — CAA 6 |      | 45 UCG  | — CGA 4  |      | stop    | 0 UAG   |      | — CUA 0 | Trp     |

|     |         |          |     |         |          |     |         |         |     |          |         |         |
|-----|---------|----------|-----|---------|----------|-----|---------|---------|-----|----------|---------|---------|
| Leu | 127 CUU | } AAG 13 | Pro | 175 CCU | } AGG 11 | His | 104 CAU | } AUG 0 | Arg | 47 CGU   | } ACG 9 |         |
|     | 187 CUC |          |     | GAG 0   |          |     | 197 CCC |         |     | GGG 0    |         | 147 CAC |
|     | 69 CUA  | — UAG 2  |     | 170 CCA | — UGG 10 |     | Gln     | 121 CAA |     | — UUG 11 | 63 CGA  | — UCG 7 |
|     | 392 CUG | — CAG 6  |     | 69 CCG  | — CGG 4  |     |         | 343 CAG |     | — CUG 21 | 115 CGG | — CCG 5 |

|     |         |          |      |         |          |     |         |          |     |         |         |
|-----|---------|----------|------|---------|----------|-----|---------|----------|-----|---------|---------|
| Ile | 165 AUU | } AAU 13 | Thr  | 131 ACU | } AGU 8  | Asn | 174 AAU | } AUU 1  | Ser | 121 AGU | } ACU 0 |
|     | 218 AUC |          |      | GAU 1   |          |     | 192 ACC |          |     | GGU 0   |         |
| Met | 71 AUA  | — UAU 5  | stop | 150 ACA | — UGU 10 | Lys | 248 AAA | — UUU 16 | Arg | 113 AGA | — UCU 5 |
|     | 221 AUG | — CAU 17 |      | 63 ACG  | — CGU 7  |     | 331 AAG | — CUU 22 |     | 110 AGG | — CCU 4 |

|     |         |          |     |         |          |     |         |         |     |          |         |         |
|-----|---------|----------|-----|---------|----------|-----|---------|---------|-----|----------|---------|---------|
| Val | 111 GUU | } AAC 20 | Ala | 185 GCU | } AGC 25 | Asp | 230 GAU | } AUC 0 | Gly | 112 GGU  | } ACC 0 |         |
|     | 146 GUC |          |     | GAC 0   |          |     | 282 GCC |         |     | GGC 0    |         | 262 GAC |
|     | 72 GUA  | — UAC 5  |     | 160 GCA | — UGC 10 |     | Glu     | 301 GAA |     | — UUC 14 | 168 GGA | — UCC 5 |
|     | 288 GUG | — CAC 19 |     | 74 GCG  | — CGC 5  |     |         | 404 GAG |     | — CUC 8  | 160 GGG | — CCC 8 |

**Figure 34** The human genetic code and associated tRNA genes. For each of the 64 codons, we show: the corresponding amino acid; the observed frequency of the codon per 10,000 codons; the codon; predicted wobble pairing to a tRNA anticodon (black lines); an unmodified tRNA anticodon sequence; and the number of tRNA genes found with this anticodon. For example, phenylalanine is encoded by UUU or UUC; UUC is seen more frequently, 203 to 171 occurrences per 10,000 total codons; both codons are expected to be decoded by a single tRNA anticodon type, GAA, using a G/U wobble; and there are 14 tRNA genes found with this anticodon. The modified anticodon sequence in the mature tRNA is not shown, even where post-transcriptional modifications can be confidently predicted (for example, when an A is used to decode a U/C third position, the A is almost certainly an inosine in the mature tRNA). The Figure also does not show the number of distinct tRNA species (such as distinct sequence families) for each anticodon; often there is more than one species for each anticodon.

# Prokaryote HMMs

- Main types of sites:
  - Codon sites
  - Translation start sites (Shine-Dalgarno)
  - Promoter elements
    - Transcription factor binding sites
  - (RNA genes / RNA folding sites)
  - (replication origin)

- Simple 7-state prokaryote genome model:
  - 1 state for intergenic regions
  - 3 states for codon positions in top-strand genes
  - 3 for codon positions in bottom-strand genes

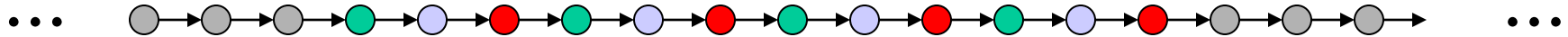
| Amino Acid | Obs/Exp | 1 <sup>st</sup> codon base | 2 <sup>nd</sup> codon base | 3 <sup>rd</sup> codon base | # codons |
|------------|---------|----------------------------|----------------------------|----------------------------|----------|
| E          | 1.92    | G                          | A                          | R                          | 2        |
| K          | 1.80    | A                          | A                          | R                          | 2        |
| D          | 1.62    | G                          | A                          | Y                          | 2        |
| M          | 1.46    | A                          | T                          | G                          | 1        |
| N          | 1.37    | A                          | A                          | Y                          | 2        |
| F          | 1.25    | T                          | T                          | Y                          | 2        |
| Q          | 1.22    | C                          | A                          | R                          | 2        |
| I          | 1.16    | A                          | T                          | Not G                      | 3        |
| A          | 1.14    | G                          | C                          | N                          | 4        |
| G          | 1.05    | G                          | G                          | N                          | 4        |
| V          | .99     | G                          | T                          | N                          | 4        |
| Y          | .98     | T                          | A                          | Y                          | 2        |
| L          | .95     | C(T)                       | T                          | N                          | 6        |
| T          | .88     | A                          | C                          | N                          | 4        |
| W          | .79     | T                          | G                          | G                          | 1        |
| P          | .74     | C                          | C                          | N                          | 4        |
| S          | .73     | T(A)                       | C(G)                       | N                          | 6        |
| H          | .67     | C                          | A                          | Y                          | 2        |
| R          | .53     | C(A)                       | G                          | N                          | 6        |
| C          | .52     | T                          | G                          | Y                          | 2        |

# Average codon biases (*lecture 3*)

- At codon position 1,
  - purines (*A* and *G*) predominate among over-represented amino acids,
  - pyrimidines (*C* and *T*) among under-represented amino acids.
- At codon position 2,
  - *A* and *T* predominate among over-represented amino acids,
  - *C* and *G* among under-represented amino acids.
- Hypotheses to explain *RWR* codon preference:
  - (Neutralist) Vestige of ancestral code? (Shepherd)
  - (Selectionist) More efficiently translated?

- These biases are somewhat subtle – but strong enough to (often) distinguish
  - coding sequences (of reasonable length)from
  - background sequence

# 7-state model for prokaryote genomes



- intergenic
- first codon position – top strand coding sequence
- second codon position – top strand coding sequence
- third codon position – top strand coding sequence
- first codon position – bottom strand coding sequence
- second codon position – bottom strand coding sequence
- third codon position – bottom strand coding sequence

a (very short!) ‘bottom-strand’ gene, in a different region of the genome:





- N.B. the emitted symbols are always *top strand* nucleotides!

# A better HMM!

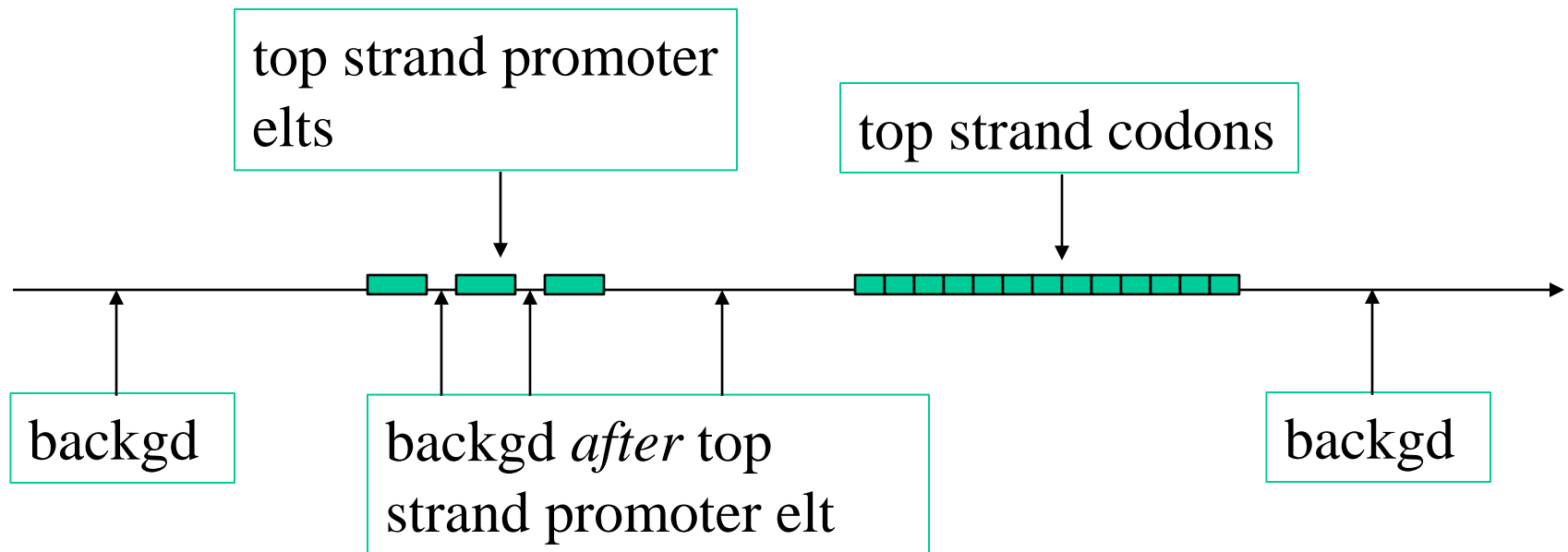
- Amino acid-specific codon blocks
  - Not really ‘sites’ as previously defined – may have more than one tRNA reader
  - Split the three 6-codon amino acids into 2 sites (4 + 2)
    - E.g. Leu: CTN and TTR ‘sites’
    - A single YTN site would also emit Phe codons
  - The other 17 aas are each 1 site
- ‘Start’ codon: NTG
  - Part of Shine-Dalgarno
- 2 Stop codon sites: TAR, TGA
- Total codon sites:  $17 + 3 \times 2 + 1 + 2 = 26$

# The Genetic Code

|   | U   | C   | A    | G    |   |
|---|-----|-----|------|------|---|
| U | Phe | Ser | Tyr  | Cys  | U |
|   | Phe | Ser | Tyr  | Cys  | C |
|   | Leu | Ser | Stop | Stop | A |
|   | Leu | Ser | Stop | Trp  | G |
| C | Leu | Pro | His  | Arg  | U |
|   | Leu | Pro | His  | Arg  | C |
|   | Leu | Pro | Gln  | Arg  | A |
|   | Leu | Pro | Gln  | Arg  | G |
| A | Ile | Thr | Asn  | Ser  | U |
|   | Ile | Thr | Asn  | Ser  | C |
|   | Ile | Thr | Lys  | Arg  | A |
|   | Met | Thr | Lys  | Arg  | G |
| G | Val | Ala | Asp  | Gly  | U |
|   | Val | Ala | Asp  | Gly  | C |
|   | Val | Ala | Glu  | Gly  | A |
|   | Val | Ala | Glu  | Gly  | G |

- Total codon *states*
- =  $26 \text{ sites} \times 2 \text{ strands} \times 3 \text{ pos} = 156$
- Transitions within & between codons are the obvious ones
  - Unless one wishes to allow for frameshift sequencing errors!
- Also, states for promoter element sites
  - TF binding sites
- Ignore RNA genes
  - (identify by sequence similarity)
- Ignore replication origins
  - Often can identify after HMM analysis, by orientation biases

- Need more than one background state, to allow memory of where one is in a gene, and strand



- May need additional backgd states if promoter element *order* is important
- Role of ‘memory’ is to *reduce* impact of biologically implausible paths
  - Model may still work without these complications – but with reduced power
- Reasonable to constrain all backgd states to have *same emission probs*

- Use Viterbi or Baum-Welch training
  - (with appropriate top vs bottom constraints etc)
- to find
  - Codon biases, aa freqs
  - Promoter elements
    - Include sites of size  $\sim 6$ , random initial emission probs
  - Shine-Dalgarno sequence preferences

# Complications in Eukaryotes

- 5' & 3' splice sites
- poly A sites
- introns
  - Must retain memory of where codon is interrupted!
- 5' & 3' UTR
- G+C variation



- Not difficult to set up an HMM with states corresponding to the above; *but* complex site spatial relationships are not well captured by Markov transition model:
  - Intron size constraints
  - Enhancers (possibly intronic!)
- Also:
  - alternative splicing
  - alternative promoters
  - overlapping sites

imply any single parse is incomplete