# HMMs for Gene Finding

**Bacterial Genome:**



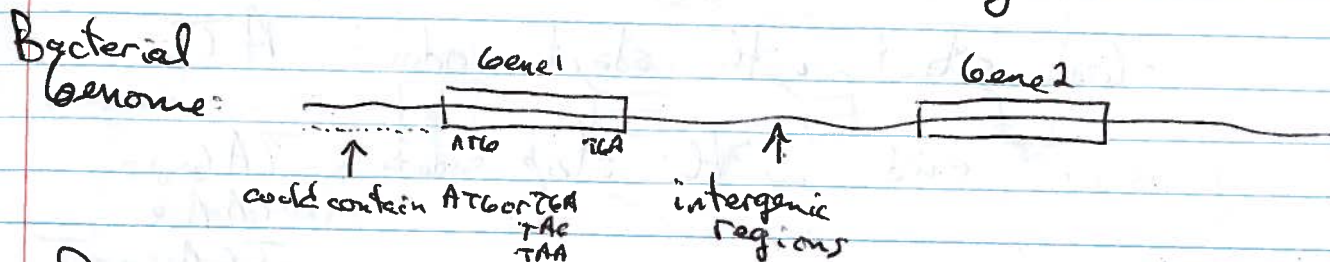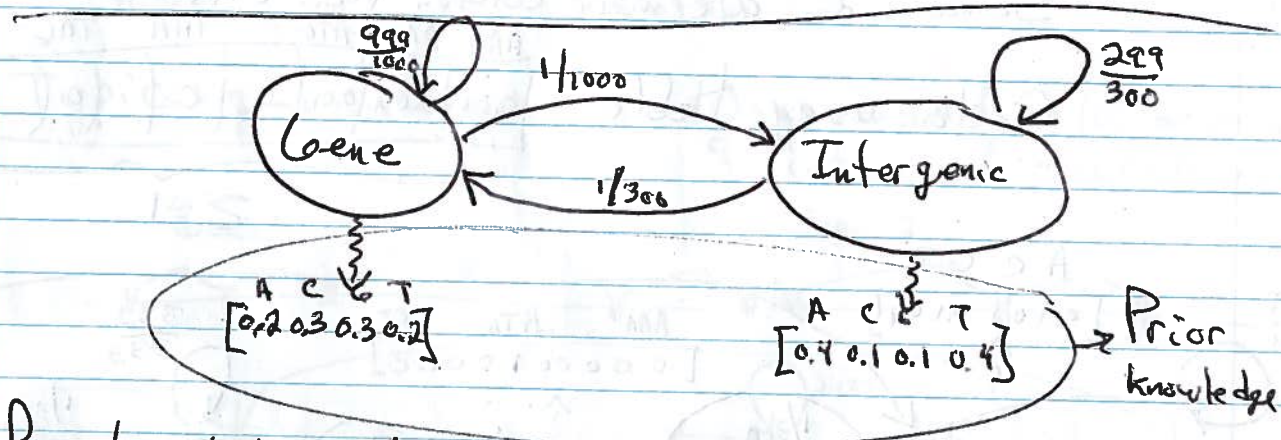could contain ATGor TGA
TAC
TAA

intergenic regions

**Problem:** **Given:** Genome sequence $X$

**Find:** Start/End position of each gene in $X$



**Prior knowledge:** Avg. gene length $\sim 1000 bp$
Avg. intergenic length $\sim 300 bp$

Genome $X = ATAGATAACA\ GGCTCG\ TGGTC\ |\ ATAT$

Viterbi Path $= IIIIITII\ |\ GGGGGGGGG\ GGG\ |\ III$

            Interg              Gene          Interg

# Gene Properties:

- Gene start with start codon: ATG

- end with stop codons: TAG
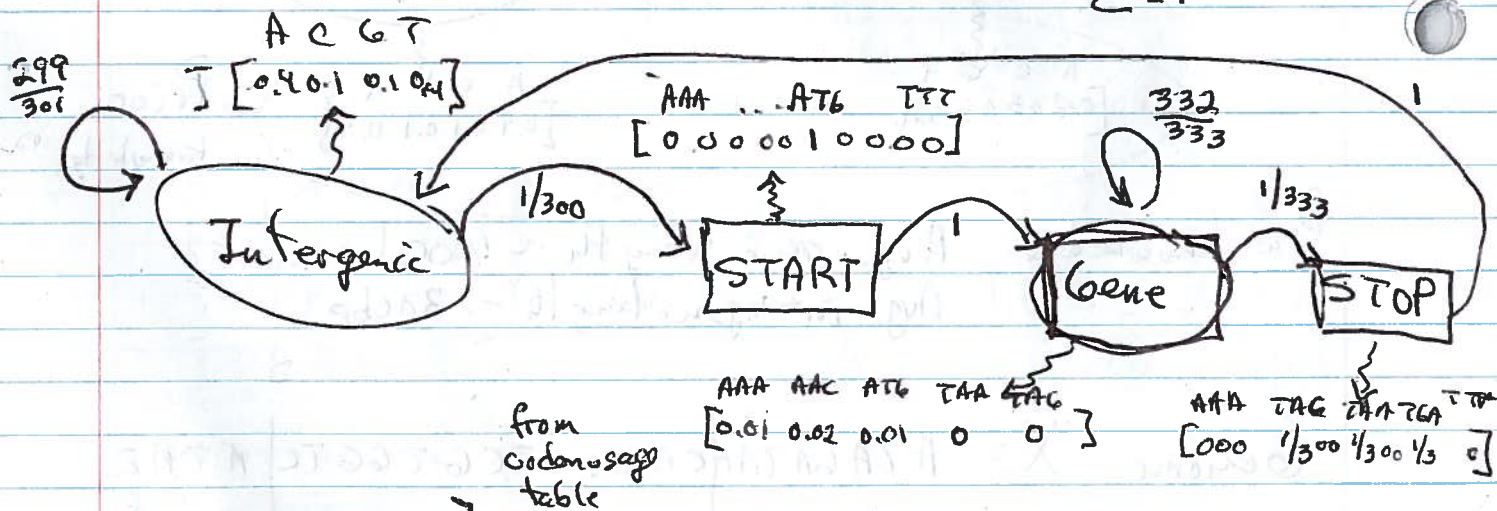  TAA
  TGA

- Gene is made of codons (triplets of nuc.)

- Some codons are more common than others

Codon usage table
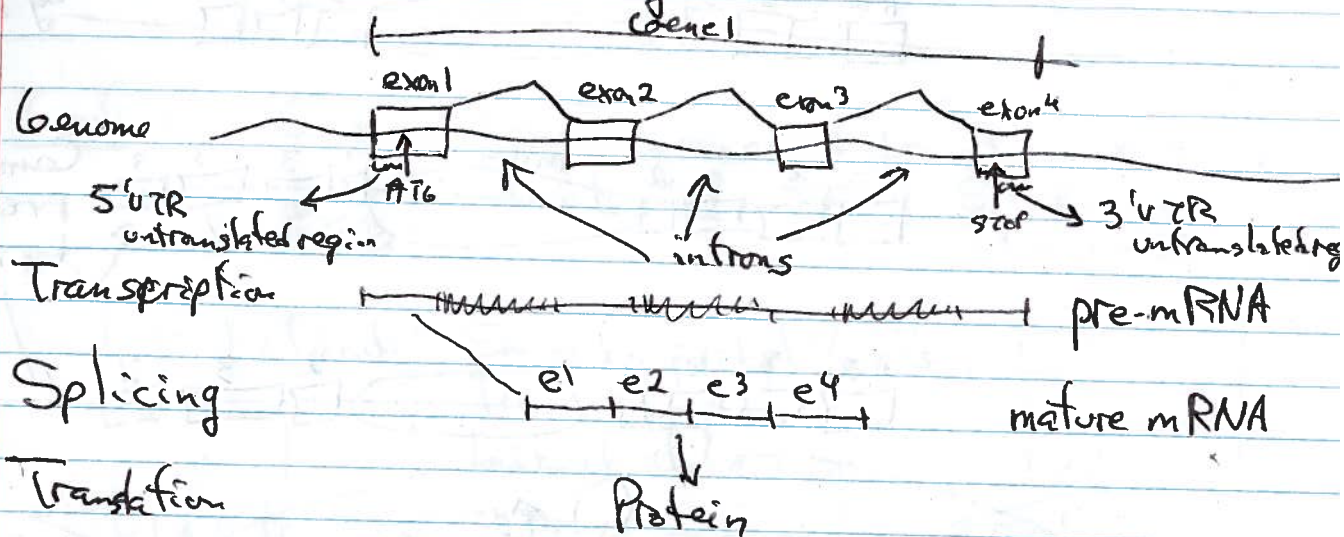
| AAA | AAC | ATG | ... | TAA | TAG | TGA |
|-----|-----|-----|-----|-----|-----|-----|
| 0.01 | 0.02 | 0.01 | | 0 | 0 | 0 |

$$\sum = 1$$



$\frac{299}{301}$

A C G T
[0.4 0.1 0.1 0.4]

AAA ... ATG  TTT
[0 0 0 0 1 0 0 0 0]

$\frac{332}{333}$   1

1/300

Intergenic

1/333

START

Gene

STOP

from codon usage table

AAA AAC ATG TAA TAG
[0.01 0.02 0.01 0 0]

AAA TAG TAA TGA TTT
[0 0 0 1/300 1/300 1/3 0]

# Gene structure in eukaryote

Genome

5'UTR untranslated region

exon1  exon2  exon3  exon4
ATG                        stop → 3'UTR untranslated region

introns

Transcription ⊢ mmmm mmmm mmmm ⊣ pre-mRNA

Splicing | e1 | e2 | e3 | e4 | mature mRNA

Translation → Protein

Typical gene in human:   10 exons of ~100 bp on average
                          9 introns of ~10000 bp on average

Gene1              Gene2

Intergenic → START → Exon → STOP

$\frac{31}{33}$

$\frac{1}{333}$

$\xi_1$

$\xi_3$

$\frac{1}{33}$  $\frac{1}{10000}$

$\xi_3$

$\xi_3$

Intron

$A C G T$
$[0.4 \ 0.1 \ 0.1 \ 0.4]$

$\frac{9999}{10000}$

ATG 3 3 3    intron    3 3 3    Can be predicted by HMM

ATG 3 3 2    intron    1 3 3 ?    Cannot be predicted by HMM

ATG 3 3 1    intron    2 3 3

I → START → EXON 3 → STOP
Exon1  Exon2    Exon1  Exon2
Intron

Problem: Could enter intron from Exon1 and leave from Exon1
codon of 2 bp

I → START → Exon 3 → STOP
Exon1  Exon2    Exon2  Exon1
Intron0   Intron1   Intron2