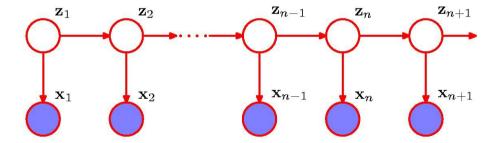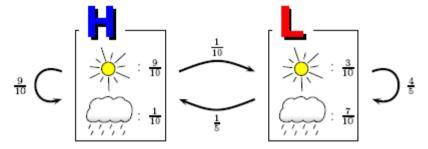# Hidden Markov Models

## Selecting the initial model parameters

## Using HMMs for gene finding

# HMMs as a generative model

A HMM *generates a sequence of observables* by moving from latent state to latent state according to the transition probabilities and *emitting an observable* (from a discrete set of observables, i.e. a finite alphabet) from each latent state visited *according to the emission probabilities* of the state ...

Model *M*:



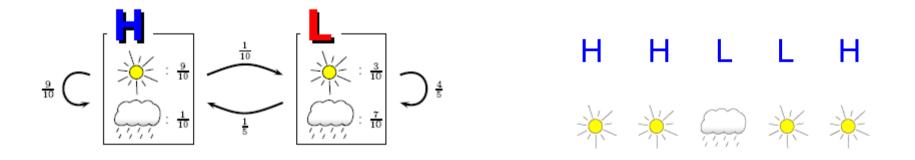A run follows a sequence of states:

H    H    L    L    H

And emits a sequence of symbols:



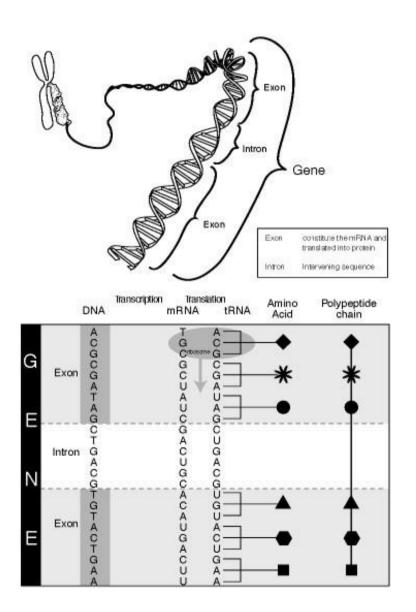For a HMM that generates finite strings (e.g. a HMM with an end-state), the language L = {**X** | $p$(**X**) > 0} is regular ...

# Selecting initial model parameters

The initial selection of transition and emission probabilities, i.e. A, π, Φ, should model (how we see) the underlying structure of the observations, i.e. the syntax of possible sequences of observations, recall that the language L = {x | P(x | θ) > 0} is regular.



The initial selection of parameters is essential just to decide which parameters are 0 (or 1), i.e. to decide which transitions of emission should never (or always) be possible ...
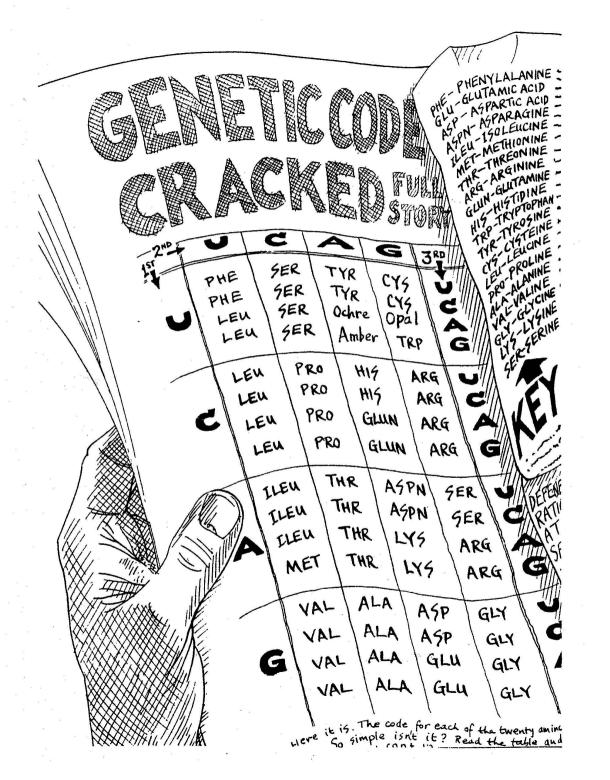
# Example – Gene finding



Each protein is encoded in a stretch of DNA. A gene ...

Which is expressed when the protein is needed ...
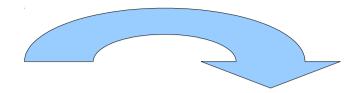
**Important problem**

Locating genes on the genome and determining how they get expressed ...

Recognizing the patterns that indicates a gene ...

>NC_002737.1 Streptococcus pyogenes M1 GAS
TTGTTGATATTCTGTTTTTTCTTTTTTAGTTTTCCACATGAAAAATAGTTGAAAACAATA
GCGGTGTCCCCTTAAAATGGCTTTTCCACAGGTTGTGGAGAACCCAAATTAACAGTGTTA
ATTTATTTTCCACAGGTTGTGGAAAAACTAACTATTATCCATCGTTCTGTGGAAAACTAG
AATAGTTTATGGTAGAATAGTTCTAGAATTATCCACAAGAAGGAACCTAGTATGACTGAA
AATGAACAAATTTTTTGGAACAGGGTCTTGGAATTAGCTCAGAGTCAATTAAAACAGGCA
ACTTATGAATTTTTTGTTCATGATGCCCGTCTATTAAAGGTCGATAAGCATATTGCAACT
ATTTACTTAGATCAAATGAAAGAGCTCTTTTGGGAAAAAAATCTTAAAGATGTTATTCTT
ACTGCTGGTTTTGAAGTTTATAACGCTCAAATTTCTGTTGACTATGTTTTCGAAGAAGAC
CTAATGATTGAGCAAAATCAGACCAAAATCAACCAAAAACCTAAGCAGCAAGCCTTAAAT
TCTTTGCCTACTGTTACTTCAGATTTAAACTCGAAATATAGTTTTGAAAACTTTATTCAA
GGAGATGAAAATCGTTGGGCTGTTGCTGCTTCAATAGCAGTAGCTAATACTCCTGGAACT
ACCTATAATCCTTTGTTTATTTGGGGTGGCCCTGGGCTTGGAAAAACCCATTTATTAAAT
GCTATTGGTAATTCTGTACTATTAGAAAATCCAAATGCTCGAATTAAATATATCACAGCT
GAAAACTTTATTAATGAGTTTGTTATCCATATTCGCCTTGATACCATGGATGAATTGAAA
GAAAAATTTCGTAATTTAGATTTACTCCTTATTGATGATATCCAATCTTTAGCTAAAAAA
ACGCTCTCTGGAACACAAGAAGAGTTCTTTAATACTTTTAATGCACTTCATAATAATAAC
AAACAAATTGTCCTAACAAGCGACCGTACACCAGATCATCTCAATGATTTAGAAGATCGA
TTAGTTACTCGTTTTAAATGGGGATTAACAGTCAATATCACACCTCCTGATTTTGAAACA
CGAGTGGCTATTTTGACAAATAAAATTCAAGAATATAACTTTATTTTTCCTCAAGATACC
ATTGAGTATTTGGCTGGTCAATTTGATTCTAATGTCAGAGATTTAGAAGGTGCCTTAAAA
GATATTAGTCTGGTTGCTAATTTCAAACAAATTGACACGATTACTGTTGACATTGCTGCC
GAAGCTATTCGCGCCAGAAAGCAAGATGGACCTAAAATGACAGTTATTCCCATCGAAGAA
ATTCAAGCGCAAGTTGGAAAATTTTACGGTGTTACCGTCAAAGAAATTAAAGCTACTAAA
CGAACACAAAATATTGTTTTAGCAAGACAAGTAGCTATGTTTTTAGCACGTGAAATGACA
GATAACAGTCTTCCTAAAATTGGAAAAGAATTTGGTGGCAGAGACCATTCAACAGTACTC
CATGCCTATAATAAAATCAAAAACATGATCAGCCAGGACGAAAGCCTTAGGATCGAAATT
GAAACCATAAAAAACAAAATTAAATAACATGTGGAAAAGAATATCTTTTATGAAATAGTT
ATCCACAAGTTGTGAACATCCATTTAGTCTTGGATTCTCTCGTTTATTTAGAGTTATCCA
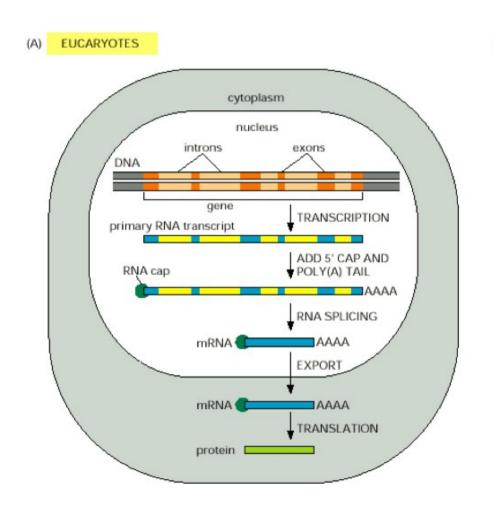CTATATACACAAGACCTACTACTACTACTTATTATTATACTTATTAAATAAAGGAGTTCT

# Viterbi decoding

>NC_002737.1 Streptococcus pyogenes M1 GAS
TTGTTGATATTCTGTTTTTTTCTTTTTTTAGTTTTCCACATGAAAAATAGTTGAAAACAATA
GCGGTGTCCCCTTAAAATGGCTTTTCCACAGGTTGTGGAGAACCCAAATTAACAGTGTTA
ATTTATTTTCCACAGGTTGTGGAAAAACTAACTATTATCCATCGTTCTGTGGAAAACTAG
AATAGTTTATGGTAGAATAGTTCTAGAATTATCCACAAGAAGGAACCTAGTATGACTGAA
AATGAACAAATTTTTTGGAACAGGGTCTTGGAATTAGCTCAGAGTCAATTAAAACAGGCA
ACTTATGAATTTTTTGTTCATGATGCCCGTCTATTAAAGGTCGATAAGCATATTGCAACT
ATTTACTTAGATCAAATGAAAGAGCTCTTTTGGGAAAAAAATCTTAAAGATGTTATTCTT
ACTGCTGGTTTTGAAGTTTATAACGCTCAAATTTCTGTTGACTATGTTTTCGAAGAAGAC
CTAATGATTGAGCAAAATCAGACCAAAATCAACCAAAAACCTAAGCAGCAAGCCTTAAAT
TCTTTGCCTACTGTTACTTCAGATTTAAACTCGAAATATAGTTTTGAAAACTTTATTCAA
GGAGATGAAAATCGTTGGGCTGTTGCTGCTTCAATAGCAGTAGCTAATACTCCTGGAACT
ACCTATAATCCTTTGTTTATTTGGGGTGGCCCTGGGCTTGGAAAAACCCATTTATTAAAT
GCTATTGGTAATTCTGTACTATTAGAAAATCCAAATGCTCGAATTAAATATATCACAGCT
GAAAACTTTATTAATGAGTTTGTTATCCATATTCGCCTTGATACCATGGATGAATTGAAA
GAAAAATTTCGTAATTTAGATTTACTCCTTATTGATGATATCCAATCTTTAGCTAAAAAA
ACGCTCTCTGGAACACAAGAAGAGTTCTTTAATACTTTTAATGCACTTCATAATAATAAC
AAACAAATTGTCCTAACAAGCGACCGTACACCAGATCATCTCAATGATTTAGAAGATCGA
TTAGTTACTCGTTTTAAATGGGGATTAACAGTCAATATCACACCTCCTGATTTTGAAACA
CGAGTGGCTATTTTGACAAATAAAATTCAAGAATATAACTTTATTTTTCCTCAAGATACC
ATTGAGTATTTGGCTGGTCAATTTGATTCTAATGTCAGAGATTTAGAAGGTGCCTTAAAA
GATATTAGTCTGGTTGCTAATTTCAAACAAATTGACACGATTACTGTTGACATTGCTGCC
GAAGCTATTCGCGCCAGAAAGCAAGATGGACCTAAAATGACAGTTATTCCCATCGAAGAA
ATTCAAGCGCAAGTTGGAAAATTTTACGGTGTTACCGTCAAAGAAATTAAAGCTACTAAA
CGAACACAAAATATTGTTTTAGCAAGACAAGTAGCTATGTTTTTTAGCACGTGAAATGACA
GATAACAGTCTTCCTAAAATTGGAAAAGAATTTGGTGGCAGAGACCATTCAACAGTACTC
CATGCCTATAATAAAATCAAAAACATGATCAGCCAGGACGAAAGCCTTAGGATCGAAATT
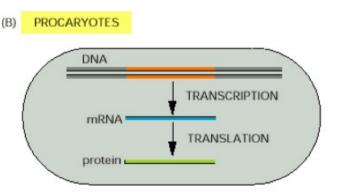GAAACCATAAAAAACAAAATTAAATAACATGTGGAAAAGAATATCTTTTATGAAATAGTT
ATCCACAAGTTGTGAACATCCATTTAGTCTTGGATTCTCTCGTTTATTTAGAGTTATCCA
CTATATACACAAGACCTACTACTACTACTTATTATTATACTTATTAAATAAAGGAGTTCT

>NC_002737.1 gene annotation Streptococcus pyogenes M1 GAS
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN

# Design a HMM that models the syntax of genes

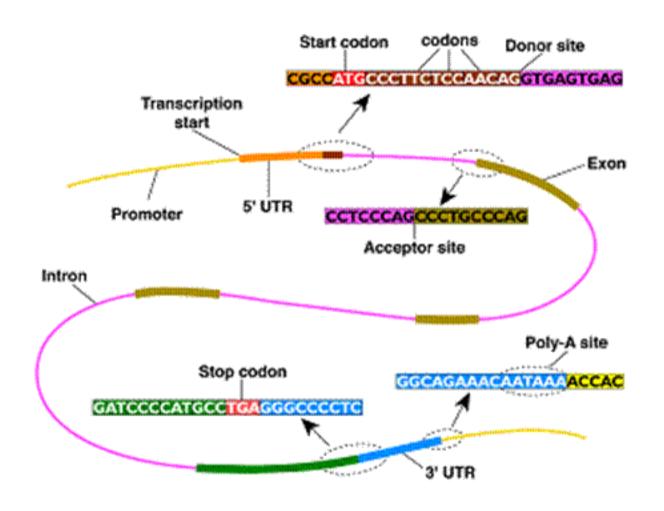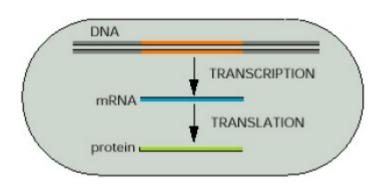# Gene structure

Depends on the organism (eucaryote or procaryote)



Large genomes. Intron/exon
structure and low coding density

Smaller genomes and
high coding density.

# Gene structure in eukaryotes



Eukaryotic gene structure in more details

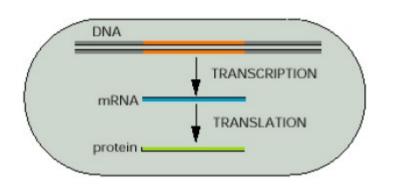# Gene structure in procaryotes

Z: NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacatgcag



C: coding

N: non-coding

A: >0
C: >0
G: >0
T: >0

A: >0
C: >0
G: >0
T: >0

$$\pi_N = 1$$
$$\pi_C = 0$$

# Gene structure in procaryotes



**Biological facts**

- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

**Z:** NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

**X:** acgatgcgctaatatgtccgatgacgtgagcataagcgacatgcag



C: coding

A: >0
C: >0
G: >0
T: >0

A: >0
C: >0
G: >0
T: >0

N: non-coding

$\pi_N = 1$
$\pi_C = 0$

# Gene structure in procaryotes



**Biological facts**

- The gene is a substring of the DNA sequence of A,C,G,T's

- The gene starts with a start-codon **atg**

**Z:** NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

**X:** acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$
$$\pi_C = 0$$



N: non-coding                    C: coding

# Gene structure in procaryotes



**Biological facts**

- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**
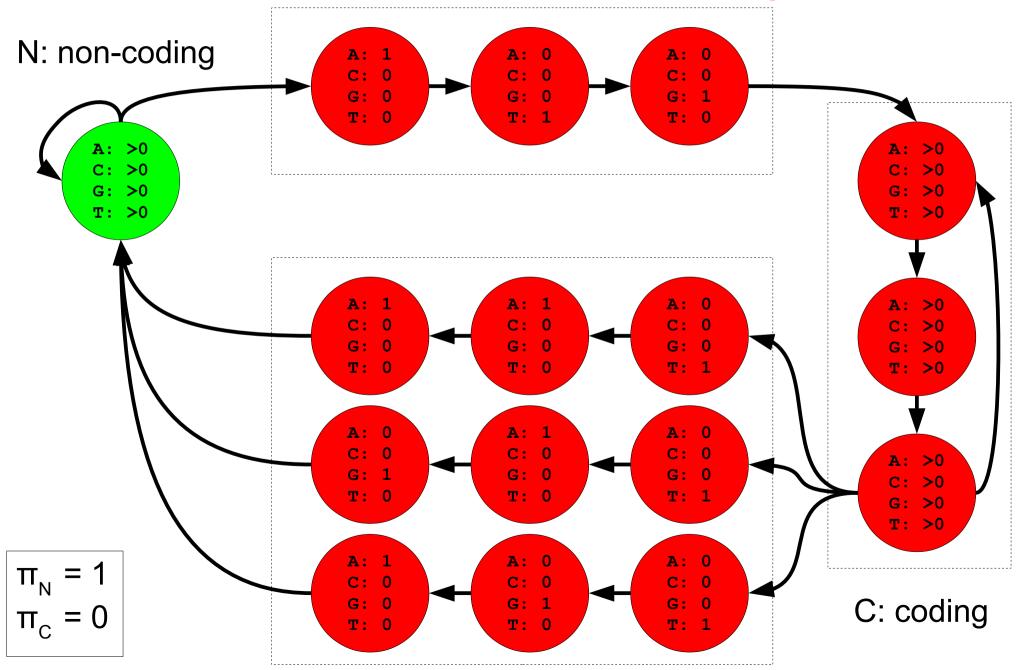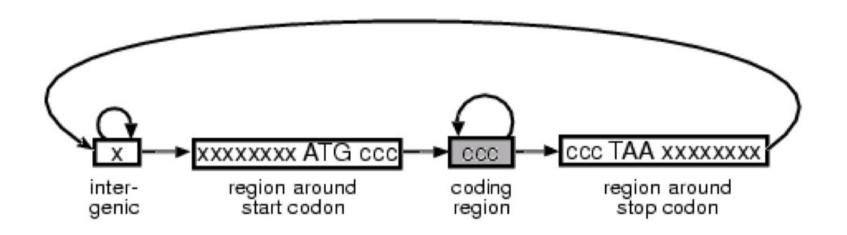- The gene ends with a stop-codon **taa**, **tag** or **tga**

Z: NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacatc

$\pi_N = 1$
$\pi_C = 0$



| A: >0 | A: 1 | A: 0 | A: 0 | A: >0 |
|-------|------|------|------|-------|
| C: >0 | C: 0 | C: 0 | C: 0 | C: >0 |
| G: >0 | G: 0 | G: 0 | G: 1 | G: >0 |
| T: >0 | T: 0 | T: 1 | T: 0 | T: >0 |

N: non-coding

C: coding

# Gene structure

N: non-coding

C: coding



$\pi_N = 1$
$\pi_C = 0$

# Gene structure

N: non-coding

C: coding



A: >0
C: >0
G: >0
T: >0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: 0
C: 0
G: 1
T: 0

A: >0
C: >0
G: >0
T: >0

A: 1
C: 0
G: 0
T: 0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: 0
C: 0
G: 1
T: 0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 1
T: 0

A: 0
C: 0
G: 0
T: 1

$\pi_N = 1$
$\pi_C = 0$

# Gene structure

- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**
- The gene ends with a stop-codon **taa**, **tag** or **tga**
- The number of nucleotides in a gene is a multiplum of 3

N: non-coding

C: coding

$\pi_N = 1$
$\pi_C = 0$

# Gene structure in procaryotes

N: non-coding

A: >0
C: >0
G: >0
T: >0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: 0
C: 0
G: 1
T: 0

A: >0
C: >0
G: >0
T: >0

A: 1
C: 0
G: 0
T: 0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: >0
C: >0
G: >0
T: >0

A: 0
C: 0
G: 1
T: 0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 0
T: 1

A: >0
C: >0
G: >0
T: >0

A: 1
C: 0
G: 0
T: 0

A: 0
C: 0
G: 1
T: 0

A: 0
C: 0
G: 0
T: 1

$\pi_N = 1$
$\pi_C = 0$

C: coding

# Gene structure in procaryotes



From "An Introduction to HMMs for Biological Sequences", A. Krogh, 1998

# Gene structure in eukaryotes



From "An Introduction to HMMs for Biological Sequences", A. Krogh, 1998

# Gene structure in procaryotes



N: non-coding
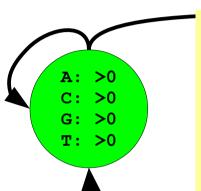
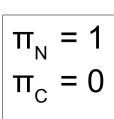C: coding

$\pi_N = 1$
$\pi_C = 0$

**Gene finding**

- Select initial model structure (e.g. as done here)

- Select model parameters by training. Either "by counting" from examples of (**X**,**Z**)'s, i.e. genes with known structure, or by EM- or Viterbi-training from examples of **X**, i.e. sequences which are known to contain a gene (as we will see later)

- Given a new sequence **X**, predict its gene structure using the Viterbi algorithm for finding the most likely sequence of underlying latent states, i.e. its gene structure

# Example – Gene finding

N: non-coding

A: 1
C: 0

A: 0
C: 0

A: 0
C: 0

A: >0
C: >0
G: >0
T: >0

A: >0
C: >0
G: >0
T: >0

A: >0
C: >0
G: >0
T: >0

A: >0
C: >0
G: >0
T: >0

$\pi_N = 1$
$\pi_C = 0$

T: 0

T: 0

T: 1

C: coding

**Gene finding**

- Select initial model structure (e.g. as done here)

- Select model parameters by training. Either "by counting" from examples of (**X**,**Z**)'s, i.e. genes with known structure, or by EM- or Viterbi-training from examples of **X**, i.e. sequences which are known to contain a gene (as we will
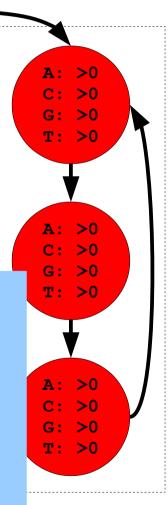
**Even more biology**

- There can be genes in both directions (and over lapping)

- There are more possible start-codons **atg**, **gtg**, and **ttg**

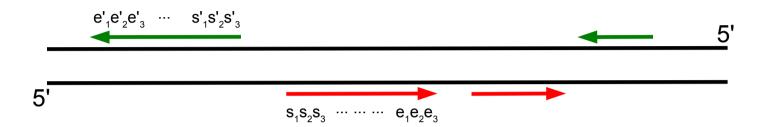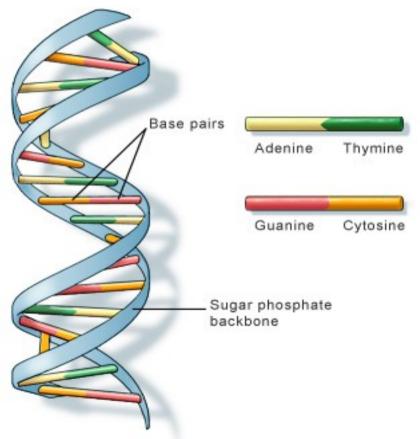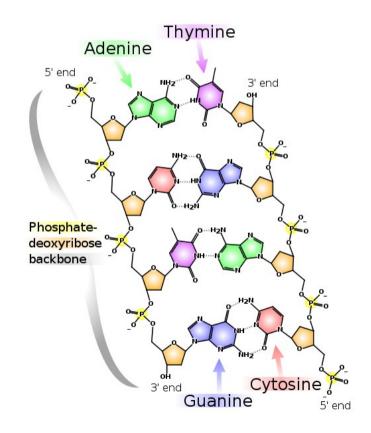- Internal codons cannot be start- or stop-codons
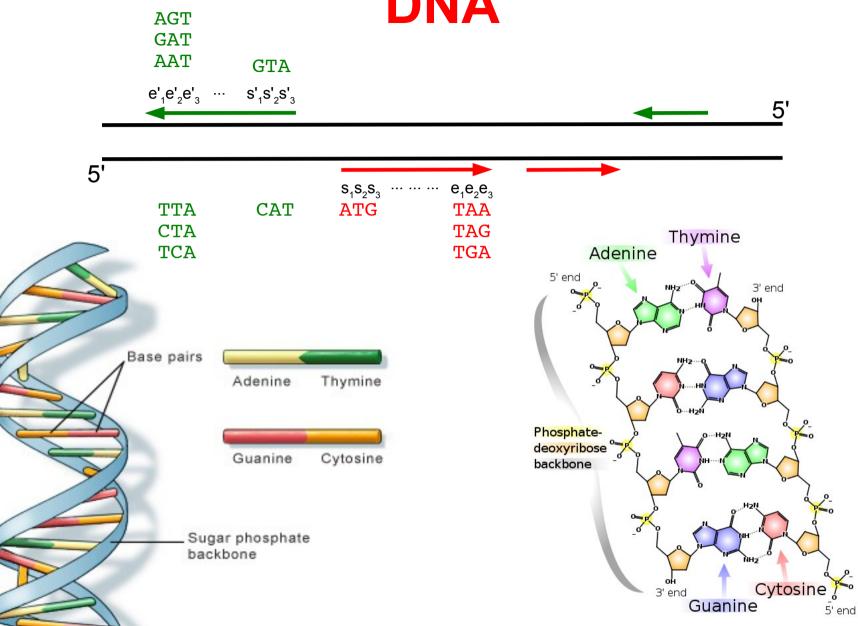
- And a lot more ...

# DNA

$e'_1 e'_2 e'_3 \cdots \quad s'_1 s'_2 s'_3$

5'

$s_1 s_2 s_3 \cdots \cdots \cdots e_1 e_2 e_3$



Base pairs

Adenine  Thymine

Guanine  Cytosine

Sugar phosphate backbone

Thymine

Adenine

5' end

3' end

Phosphate-deoxyribose backbone

3' end

Guanine

Cytosine

5' end

U.S. National Library of Medicine

# DNA

AGT
GAT
AAT     GTA

$e'_1 e'_2 e'_3$   ⋯   $s'_1 s'_2 s'_3$

5'

5'

$s_1 s_2 s_3$ ⋯ ⋯ ⋯ $e_1 e_2 e_3$

TTA     CAT     ATG        TAA
CTA                           TAG
TCA                           TGA



Base pairs

Adenine    Thymine

Guanine    Cytosine

Sugar phosphate
backbone



Thymine

Adenine

5' end

3' end

Phosphate-
deoxyribose
backbone

3' end

Guanine

Cytosine

5' end

U.S. National Library of Medicine

C: coding left-to-right

**Even more biology**

There can be genes in both directions

N: Non-coding

R: coding right-to-left

$\pi_N = 1$

$\pi_C = 0$
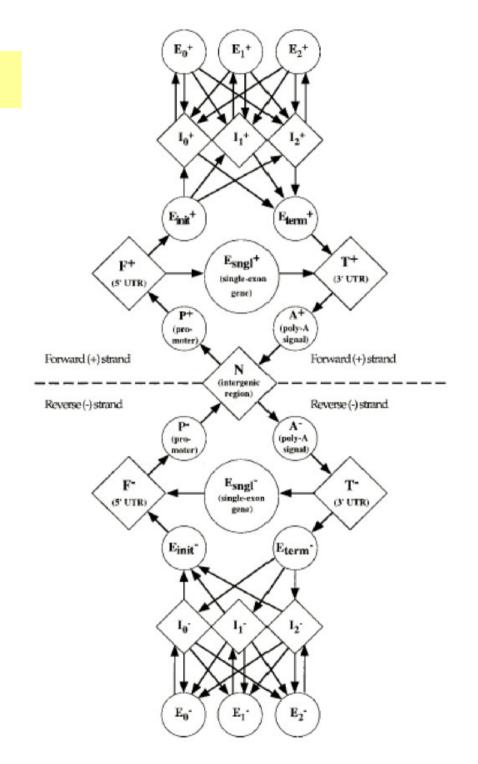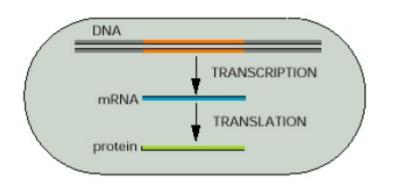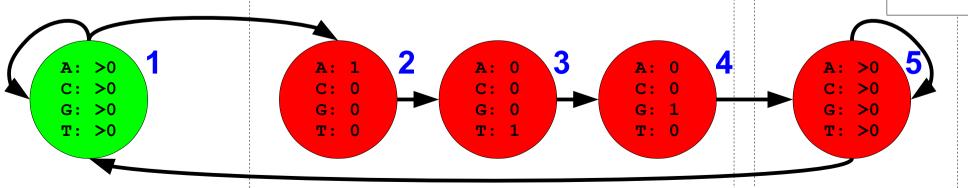
Eukaryotic gene structure in both directions (GenScan)

# Problem: From annotation to Z



## Biological facts

- The gene is a substring of the DNA sequence of A,C,G,T's

- The gene starts with a start-codon **atg**

**Z:**  NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCNNNNNNNNNNNN

**X:**  acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$
$$\pi_C = 0$$



A: >0  C: >0  G: >0  T: >0    **1**

A: 1  C: 0  G: 0  T: 0    **2**

A: 0  C: 0  G: 0  T: 1    **3**

A: 0  C: 0  G: 1  T: 0    **4**

A: >0  C: >0  G: >0  T: >0    **5**

N: non-coding

C: coding
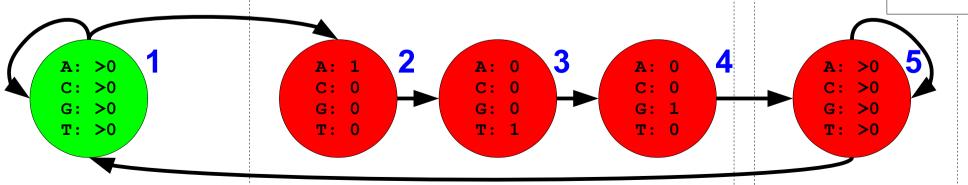
# Problem: From annotation to Z

**Problem:** The string **Z=NNNCCC....** is not a prober sequence of states in the illustrated HMM, but is can easily be converted into one.

protein

**Z:** NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCNNNNNNNNNNNN

**X:** acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$\pi_N = 1$
$\pi_C = 0$



**1**
A: >0
C: >0
G: >0
T: >0

**2**
A: 1
C: 0
G: 0
T: 0

**3**
A: 0
C: 0
G: 0
T: 1

**4**
A: 0
C: 0
G: 1
T: 0

**5**
A: >0
C: >0
G: >0
T: >0

N: non-coding

C: coding

# Problem: From annotation to Z

**Problem:** The string **Z**=NNNCCC.... is not a prober sequence of states in the illustrated HMM, but is can easily be converted into one.

protein

```
111234555555111111112345555555555511111111111
```

Z : NNNCCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCNNNNNNNNNNN

X : acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$\pi_N = 1$
$\pi_C = 0$



1
A: >0
C: >0
G: >0
T: >0

2
A: 1
C: 0
G: 0
T: 0

3
A: 0
C: 0
G: 0
T: 1

4
A: 0
C: 0
G: 1
T: 0

5
A: >0
C: >0
G: >0
T: >0

N: non-coding

C: coding

# Evaluating performance



Evaluation of Gene Structure Prediction Programs (Burset and Guigo, 1996)