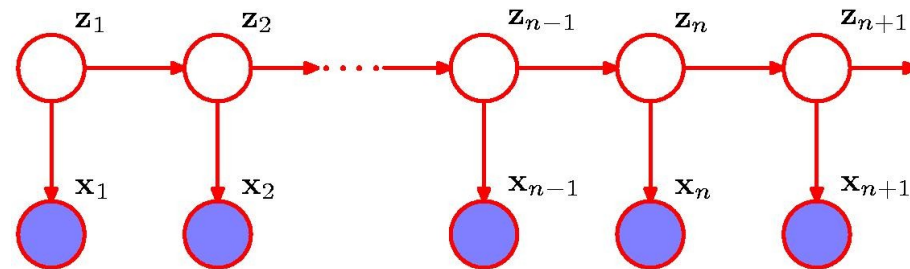


Hidden Markov Models

Selecting the initial model parameters

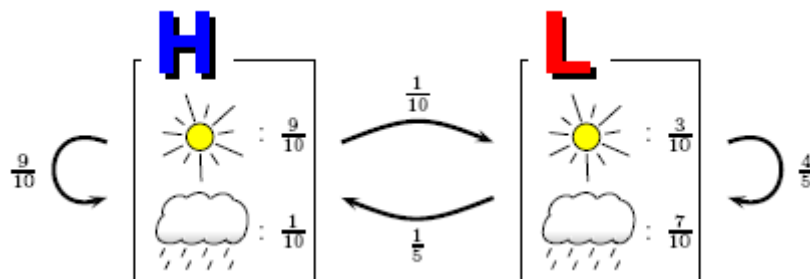
Using HMMs for (simple) gene finding



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

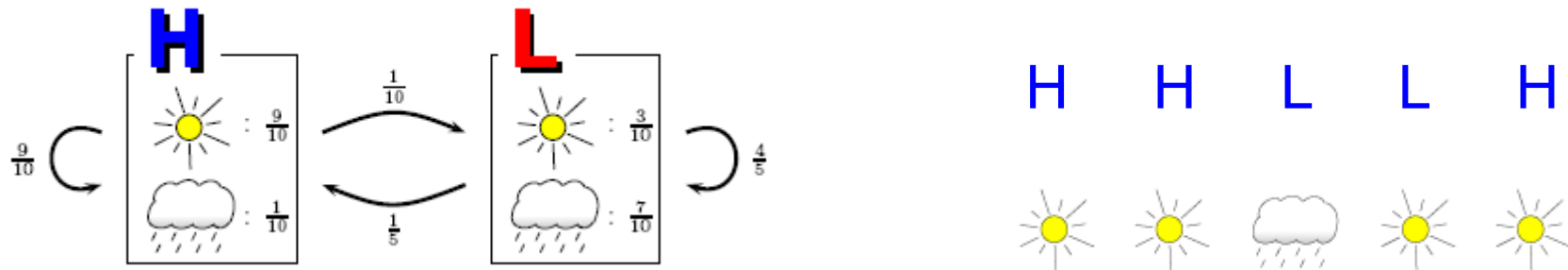
And emits a sequence of symbols:



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

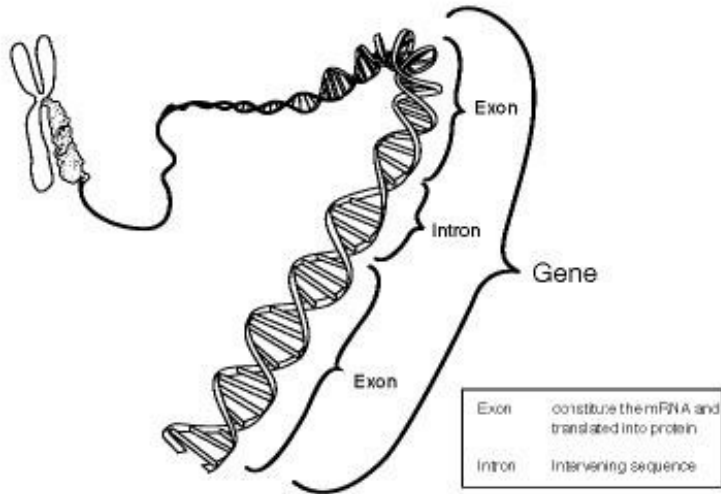
Selecting initial model parameters

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 $L = \{x \mid P(x \mid \theta) > 0\}$ is regular.



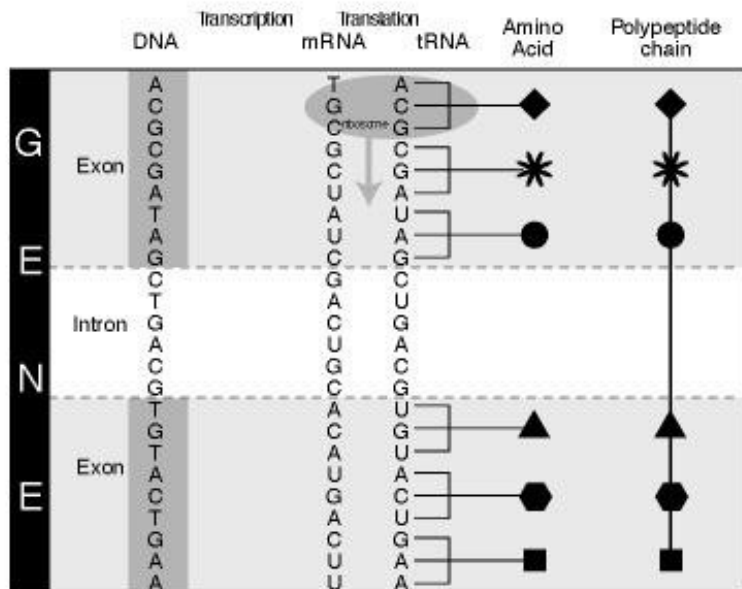
The initial selection of parameters is essential just to decide which parameters are 0 (or 1), i.e. to decide which transitions or emissions should never (or always) be possible ...

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 ...

GENETIC CODE CRACKED FULL STORY

2ND 1ST	U	C	A	G	3RD
U	PHE PHE LEU LEU	SER SER SER SER	TYR TYR Ochre Amber	CYS CYS Opal TRP	U C A G
C	LEU LEU LEU LEU	PRO PRO PRO PRO	HIS HIS GLUN GLUN	ARG ARG ARG ARG	U C A G
A	ILEU ILEU ILEU MET	THR THR THR THR	ASPN ASPN LYS LYS	SER SER ARG ARG	U C A G
G	VAL VAL VAL VAL	ALA ALA ALA ALA	ASP ASP GLU GLU	GLY GLY GLY GLY	U C A G

PHE - PHENYLALANINE
 GLU - GLUTAMIC ACID
 ASP - ASPARTIC ACID
 ASPN - ASPARAGINE
 ILEU - ISOLEUCINE
 MET - METHIONINE
 THR - THREONINE
 ARG - ARGinine
 GLUN - GLUTAMINE
 HIS - HISTIDINE
 TRP - TRYPTOPHAN
 TYR - TYROSINE
 CYS - CYSTEINE
 LEU - LEUCINE
 PRO - PROLINE
 ALA - ALANINE
 VAL - VALINE
 GLY - GLYCINE
 LYS - LYSINE
 SER - SERINE

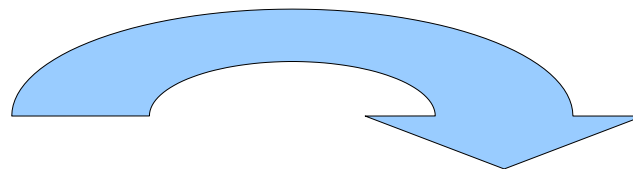
KEY

DEFENSE
RATH
AT
SF

Here it is. The code for each of the twenty amino acids. So simple isn't it? Read the table and you can't miss it.

>NC_002737.1 Streptococcus pyogenes M1 GAS
TTGTTGATATTCTGTTTTTCTTTTTAGTTTTCCACATGAAAAATAGTTGAAAACAATA
GCGGTGTCCTTAAAATGGCTTTTCCACAGGTTGTGGAGAACCCAAATTAACAGTGTTA
ATTTATTTTCCACAGGTTGTGGAAAACTAACTATTATCCATCGTTCTGTGGAAAACTAG
AATAGTTTATGGTAGAATAGTTCTAGAATTATCCACAAGAAGGAACCTAGTATGACTGAA
AATGAACAAATTTTTTGGAACAGGGTCTTGGAAATTAGCTCAGAGTCAATTAACAGGCA
ACTTATGAATTTTTGTTTCATGATGCCCGTCTATTAAGGTCGATAAGCATATTGCAACT
ATTTACTTAGATCAAATGAAAGAGCTCTTTGGGAAAAAATCTTAAAGATGTTATTCTT
ACTGCTGGTTTTGAAGTTTATAACGCTCAAATTTCTGTTGACTATGTTTTCGAAGAAGAC
CTAATGATTGAGCAAATCAGACCAAATCAACCAAACCTAAGCAGCAAGCCTTAAAT
TCTTTGCCTACTGTTACTTCAGATTTAACTCGAAATATAGTTTTGAAAACTTTATTCAA
GGAGATGAAAATCGTTGGGCTGTTGCTGCTTCAATAGCAGTAGCTAATACTCCTGGA
ACCTATAATCCTTTGTTTATTTGGGGTGGCCCTGGGCTTGGAAAAACCCATTTATTAAT
GCTATTGGTAATTCTGTACTATTAGAAAATCCAAATGCTCGAATTAATATATCACAGCT
GAAAACTTTATTAATGAGTTTGTTATCCATATTCGCCTTGATACCATGGATGAATTGAAA
GAAAAATTTGTAATTTAGATTTACTCCTTATTGATGATATCCAATCTTTAGCTAAAAAA
ACGCTCTCTGGAACACAAGAAGAGTTCTTTAATACTTTTAAATGCACTTCATAATAAAC
AAACAAATTGTCCTAACAAGCGACCGTACACCAGATCATCTCAATGATTTAGAAGATCGA
TTAGTTACTCGTTTTAAATGGGGATTAACAGTCAATATCACACCTCCTGATTTTGAAACA
CGAGTGGCTATTTTGACAAATAAAATCAAGAATATAACTTTATTTTTTCTCAAGATACC
ATTGAGTATTTGGCTGGTCAATTTGATTCTAATGTCAGAGATTTAGAAGGTGCCTTAAAA
GATATTAGTCTGGTTGCTAATTTCAAACAAATTGACACGATTACTGTTGACATTGCTGCC
GAAGCTATTCGCGCCAGAAAGCAAGATGGACCTAAAATGACAGTTATTCATCGAAGAA
ATTCAAGCGCAAGTTGGAAAATTTACGGTGTACCGTCAAAGAAATTAAGCTACTAAA
CGAACACAAAATATTGTTTTAGCAAGACAAGTAGCTATGTTTTTAGCACGTGAAATGACA
GATAACAGTCTTCCTAAAATTTGGAAAAGAATTTGGTGGCAGAGACCATTCAACAGTACTC
CATGCCTATAATAAAATCAAAAACATGATCAGCCAGGACGAAAGCCTTAGGATCGAAAT
GAAACCATAAAAAACAAAATTAATAACATGTGGAAAAGAATATCTTTTATGAAATAGTT
ATCCACAAGTTGTGAACATCCATTTAGTCTTGGATTCTCTCGTTTATTTAGAGTTATCCA
CTATATACACAAGACCTACTACTACTATTATTATACTTATTAATAAAGGAGTTCT

Viterbi decoding



```
>NC_002737.1 Streptococcus pyogenes M1 GAS
TTGTTGATATTCTGTTTTTTCTTTTTAGTTTTCCACATGAAAATAGTTGAAAACAATA
GCGGTGTCCTTAAAATGGCTTTTCCACAGGTTGTGGAGAACCCAAATTAACAGTGTTA
ATTTATTTTCCACAGGTTGTGGAAAACTAACTATTATCCATCGTTCTGTGAAAACTAG
AATAGTTTATGGTAGAATAGTTCTAGAAATTAACACAAGAAGGAACCTAGTATGACTGAA
AATGAACAAATTTTTGGAACAGGGTCTTGGAAATTAGCTCAGAGTCAATTAACAGGCA
ACTTATGAATTTTTGTTTCATGATGCCGTCTATTAAGGTCGATAAGCATATTGCAACT
ATTTACTTAGATCAAATGAAAGAGCTCTTTGGGAAAAAATCTTAAGATGTTATTCTT
ACTGCTGGTTTTGAAGTTTATAACGCTCAAATTTCTGTTGACTATGTTTTCGAAGAAGAC
CTAATGATTGAGCAAATCAGACAAAATCAACAAAAACCTAAGCAGCAAGCCTTAAAT
TCTTGCCTACTGTTACTTCAGATTTAACTCGAAATATAGTTTTGAAACTTTATTCAA
GGAGATGAAAATCGTTGGGCTGTTGCTGCTTCAATAGCAGTAGCTAATACTCCTGGAAct
ACCTATAATCCTTTGTTTATTTGGGGTGGCCCTGGGCTTGAAAAAACCCATTTATTAAT
GCTATTGGTAATTCTGACTATTAGAAAATCCAAATGCTCGAATTAATATATCACAGCT
GAAAACCTTTATTAATGAGTTTGTTATCCATATTGCCTTGATACCATGGATGAATTGAAA
GAAAAATTCGTAATTTAGATTTACTCCTTATTGATGATATCCAATCTTTAGCTAAAAAA
ACGCTCTCTGGAACACAAGAAGAGTTCTTTAATACTTTTAAATGCACTTCATAATAAAC
AAACAAATTGTCCTAACAAGCGACCGTACACCAGATCATCTCAATGATTTAGAAGATCGA
TTAGTTACTCGTTTTAAATGGGGATTAACAGTCAATATCACACCTCCTGATTTTGAAACA
CGAGTGGCTATTTTGACAAATAAAATTAAGAATATAAATTTATTTTCCCTCAAGATACC
ATTGAGTATTTGGCTGGTCAATTTGATTCTAATGTCTAGAGATTTAGAAGGTGCCTTAAAA
GATATTAGTCTGGTTGCTAATTTCAAACAAATTGACACGATTACTGTTGACATTGCTGCC
GAAGCTATTCGCGCCAGAAAGCAAGATGGACCTAAAATGACAGTTATTCCATCGAAGAA
ATTCAAGCGCAAGTTGGAAAATTTACGGTGTTACCGTCAAAGAAATTAAGCTACTAAA
CGAACACAAAATATTGTTTTAGCAAGACAAGTAGCTATGTTTTTAGCACGTGAAATGACA
GATAACAGTCTTCTAAAATTTGGAAAAGAATTTGGTGGCAGAGACCATTCAACAGTACTC
CATGCCTATAATAAAATCAAAAACATGATCAGCCAGGACGAAAGCCTTAGGATCGAAAT
GAAACCATAAAAAACAAAATTAATAACATGTGGAAAAGAATATCTTTTATGAAATAGTT
ATCCACAAGTTGTGAACATCCATTTAGTCTTGGATTCTCTCGTTTATTTAGAGTTATCCA
CTATATACACAAGACCTACTACTACTATTATTATACTTATTAATAAAGGAGTTCT
```

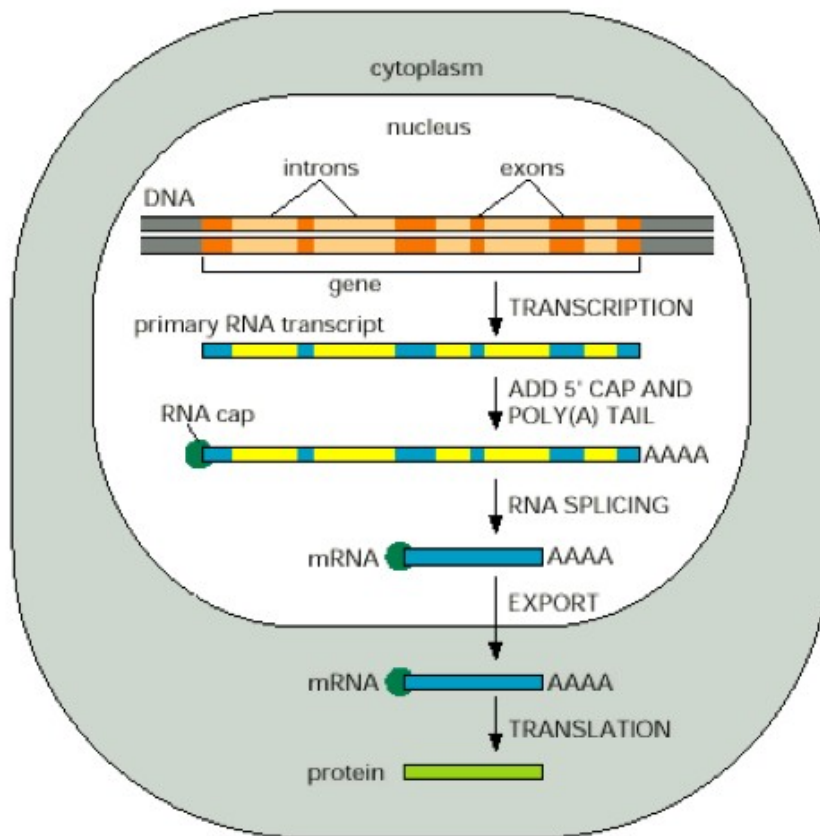
```
>NC_002737.1 gene annotation Streptococcus pyogenes M1 GAS
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
```

Design a HMM that models the syntax of genes

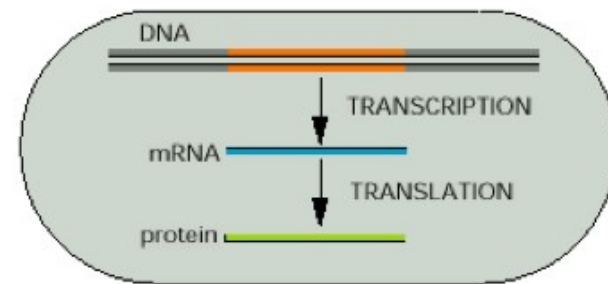
Gene structure

Depends on the organism (eucaryote or procaryote)

(A) EUCARYOTES



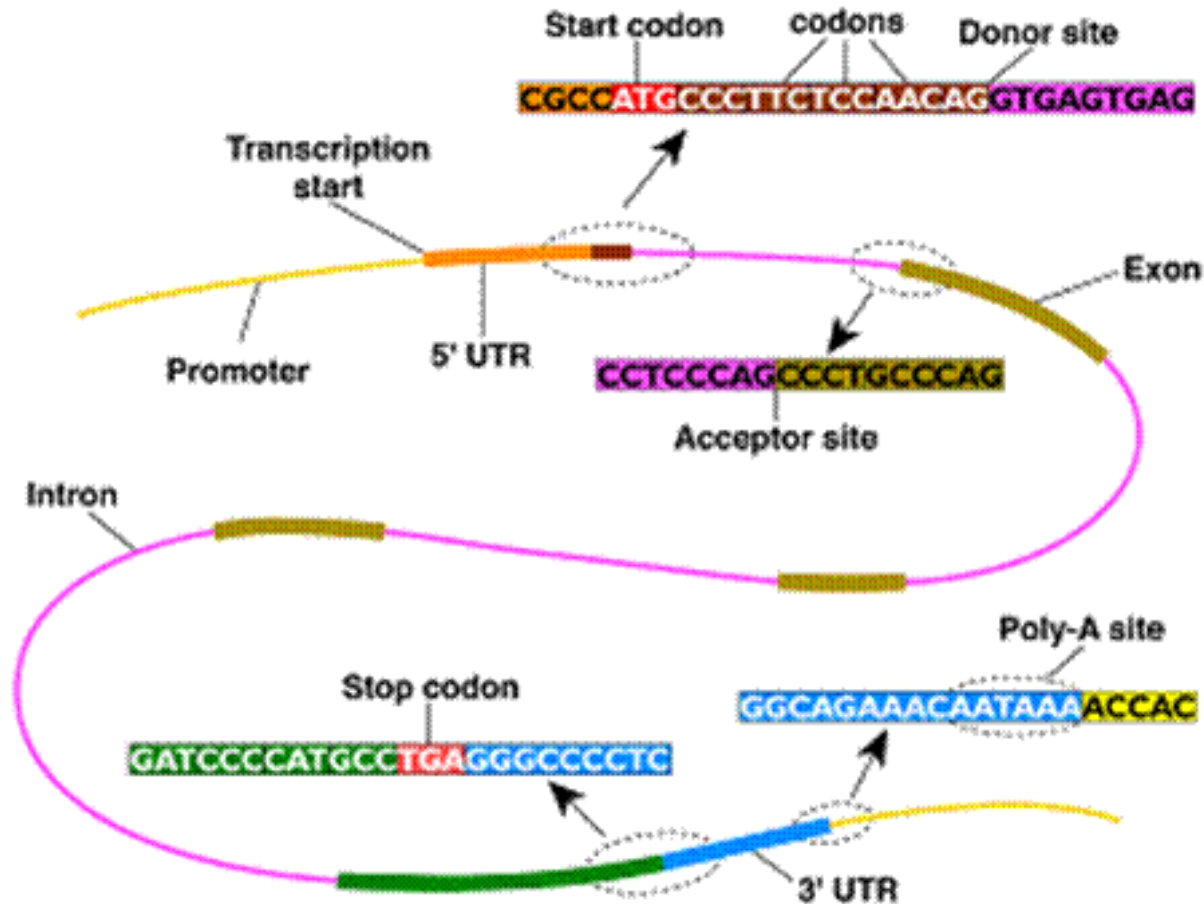
(B) PROCARYOTES



Smaller genomes and high coding density.

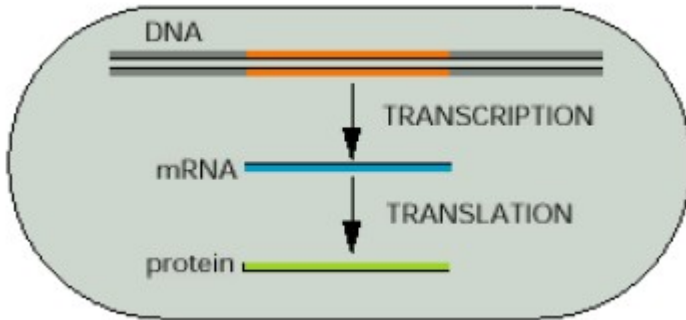
Large genomes. Intron/exon structure and low coding density

Gene structure in eukaryotes



Eukaryotic gene structure in more details

Gene structure in procaryotes

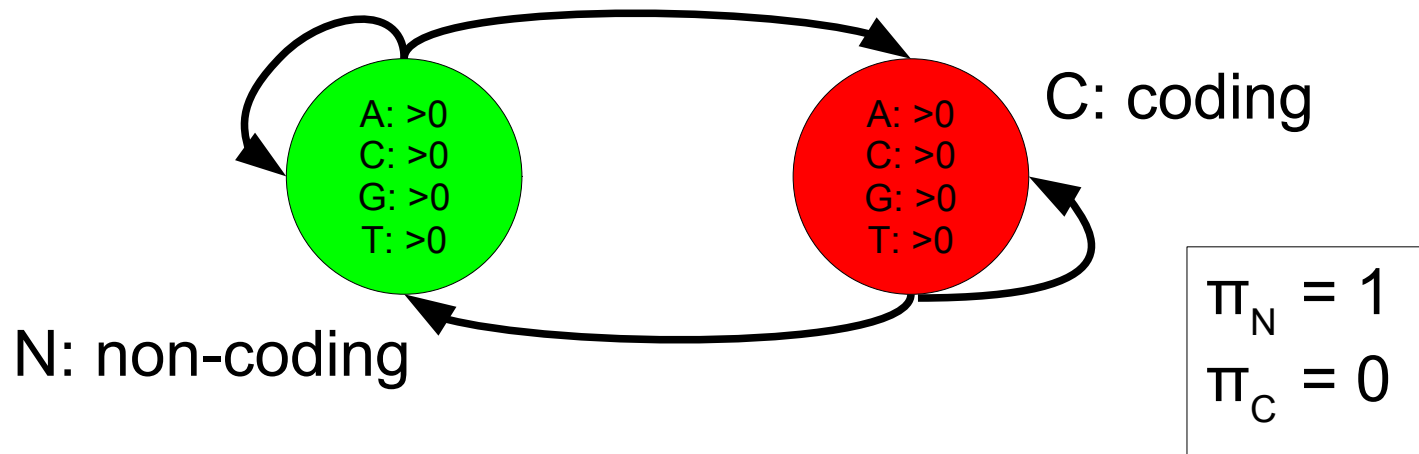


Biological facts

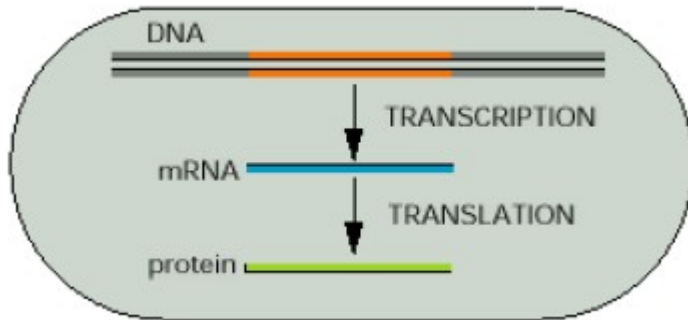
- The gene is a substring of the DNA sequence of A,C,G,T's

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacatgcag



Gene structure in procaryotes

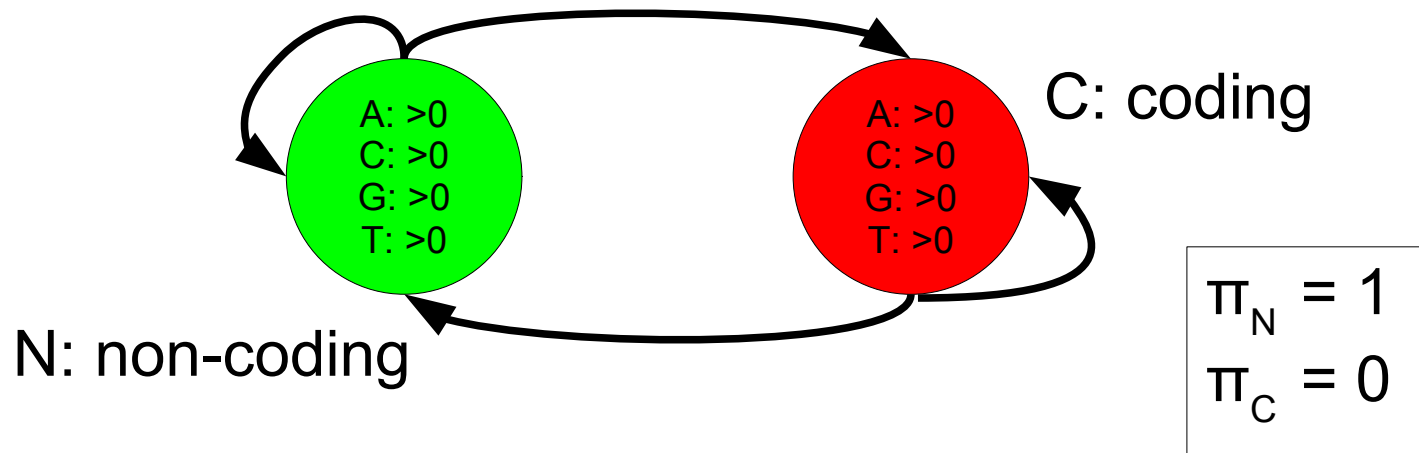


Biological facts

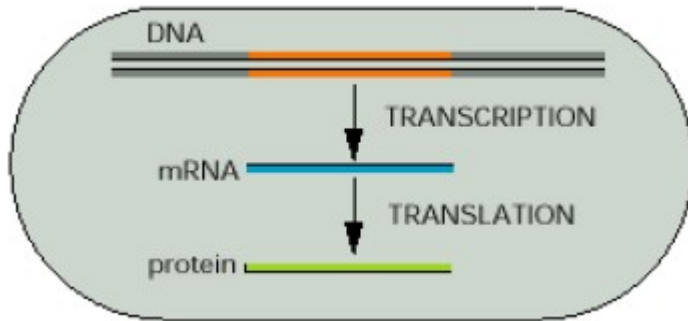
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacatgcag



Gene structure in procaryotes



Biological facts

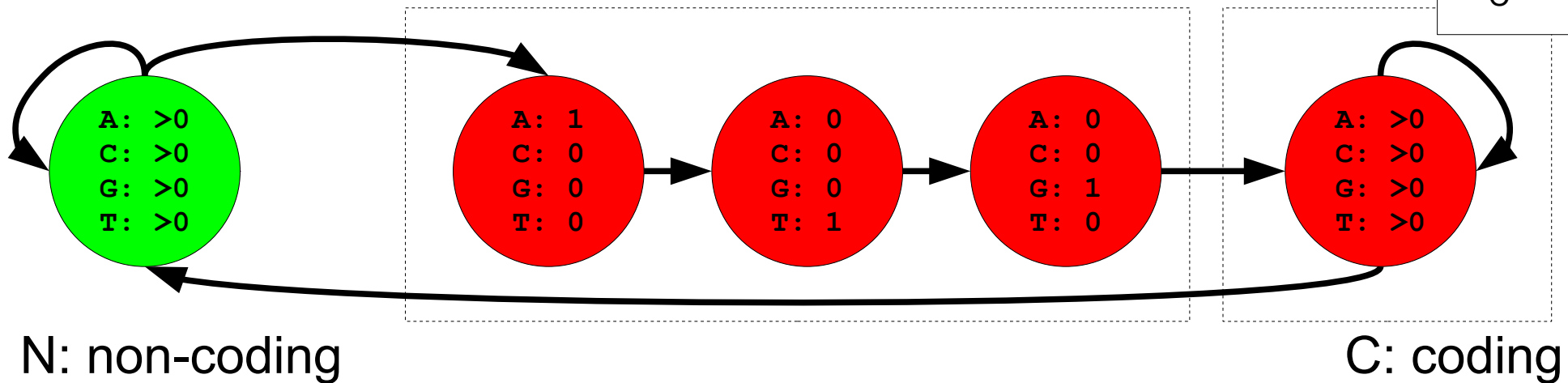
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCNNNNNNNNNN

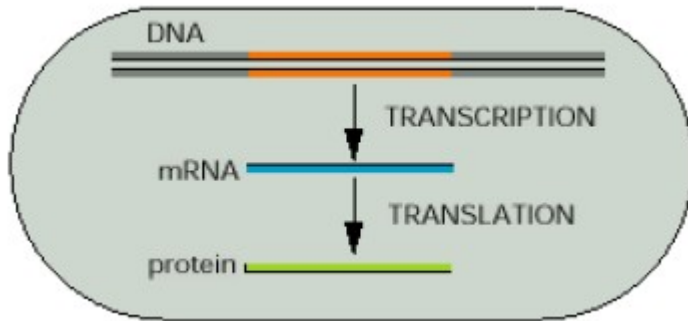
X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$



Gene structure in procaryotes



Biological facts

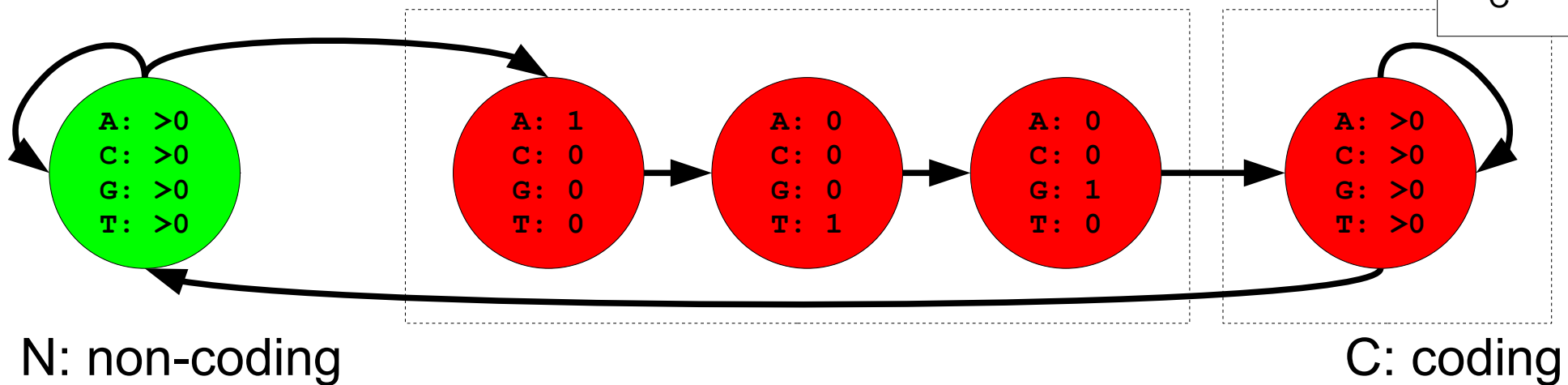
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**
- The gene ends with a stop-codon **taa, tag or tga**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$

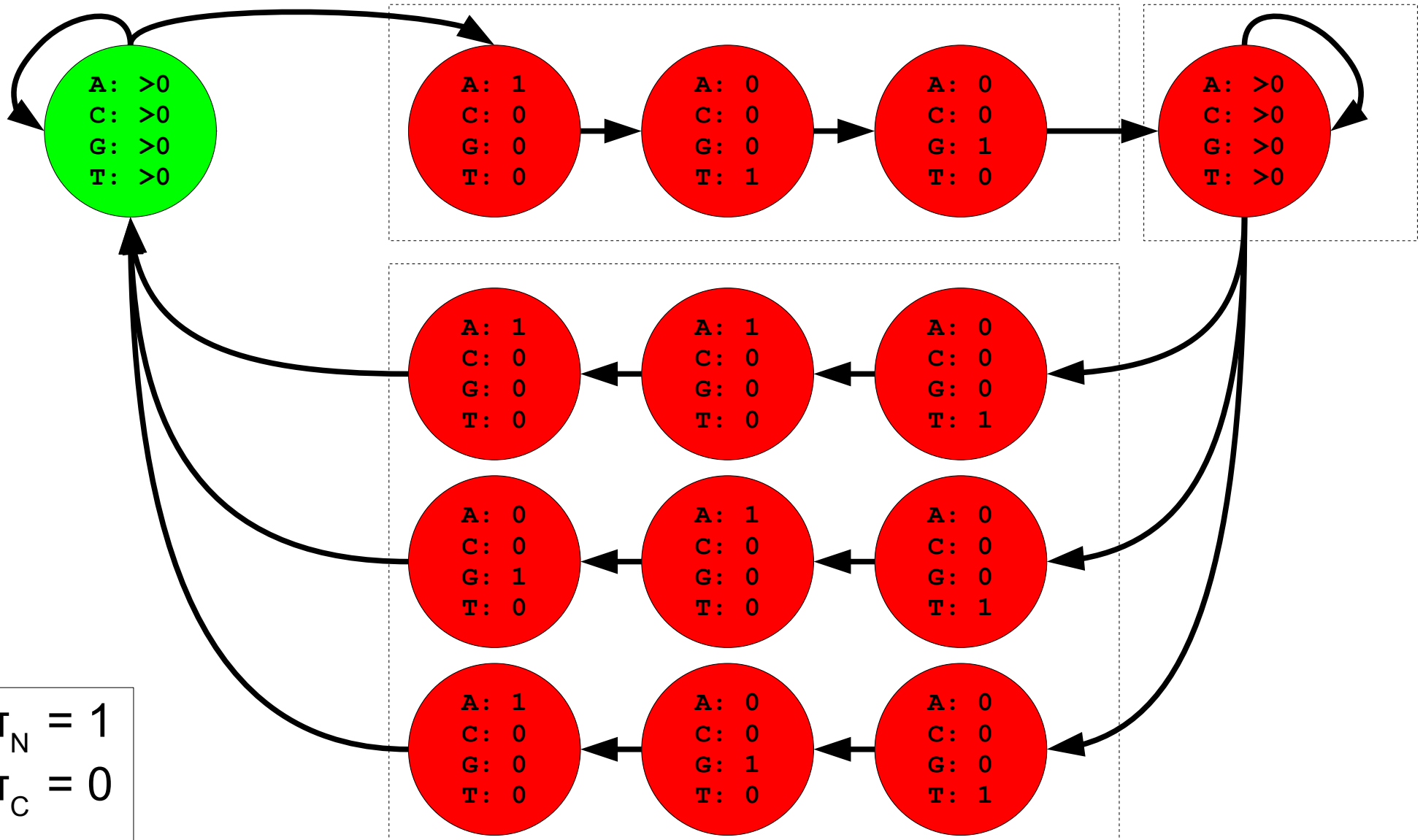


Gene structure

- 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**

N: non-coding

C: coding

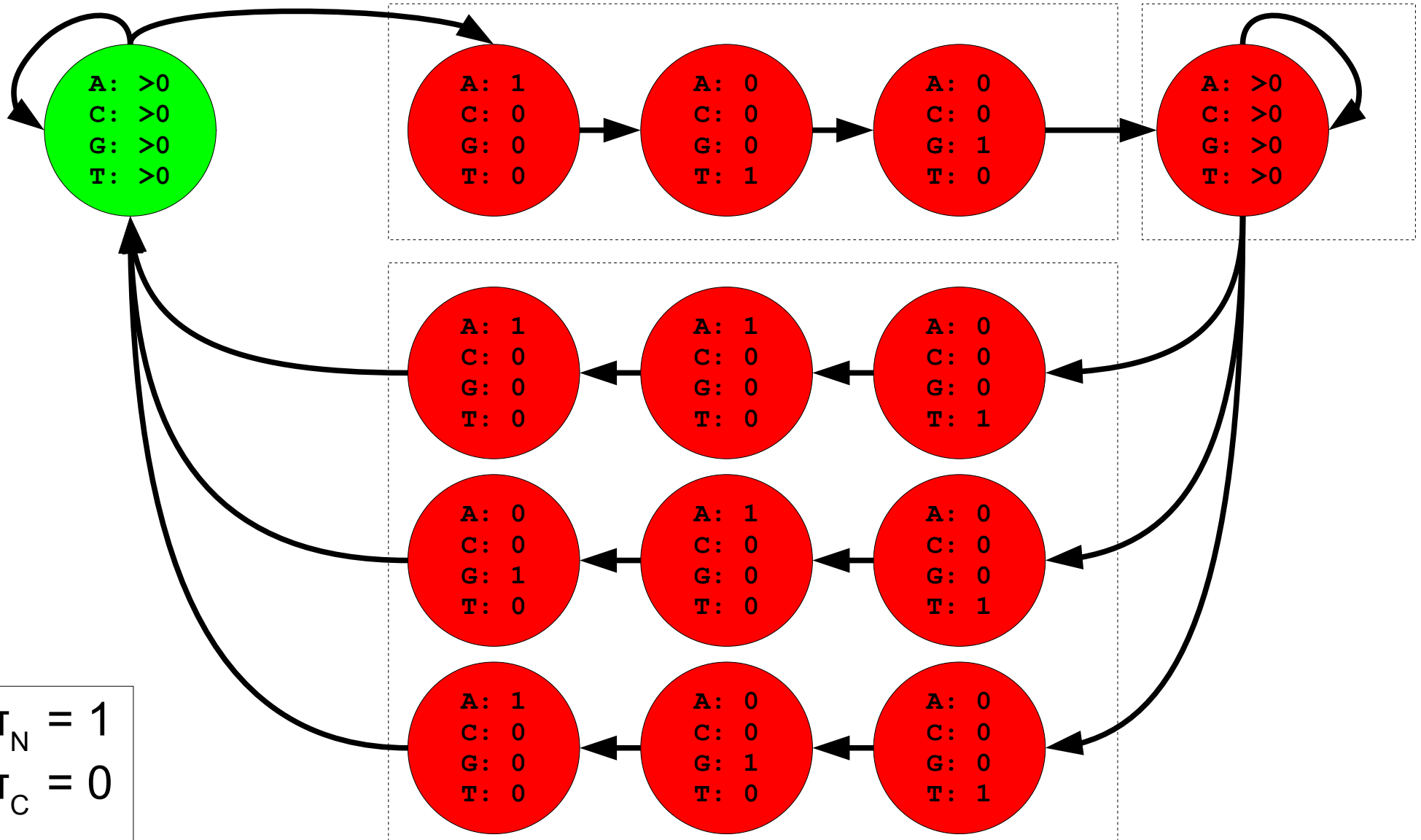


Gene structure

- 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 multiple of 3

N: non-coding

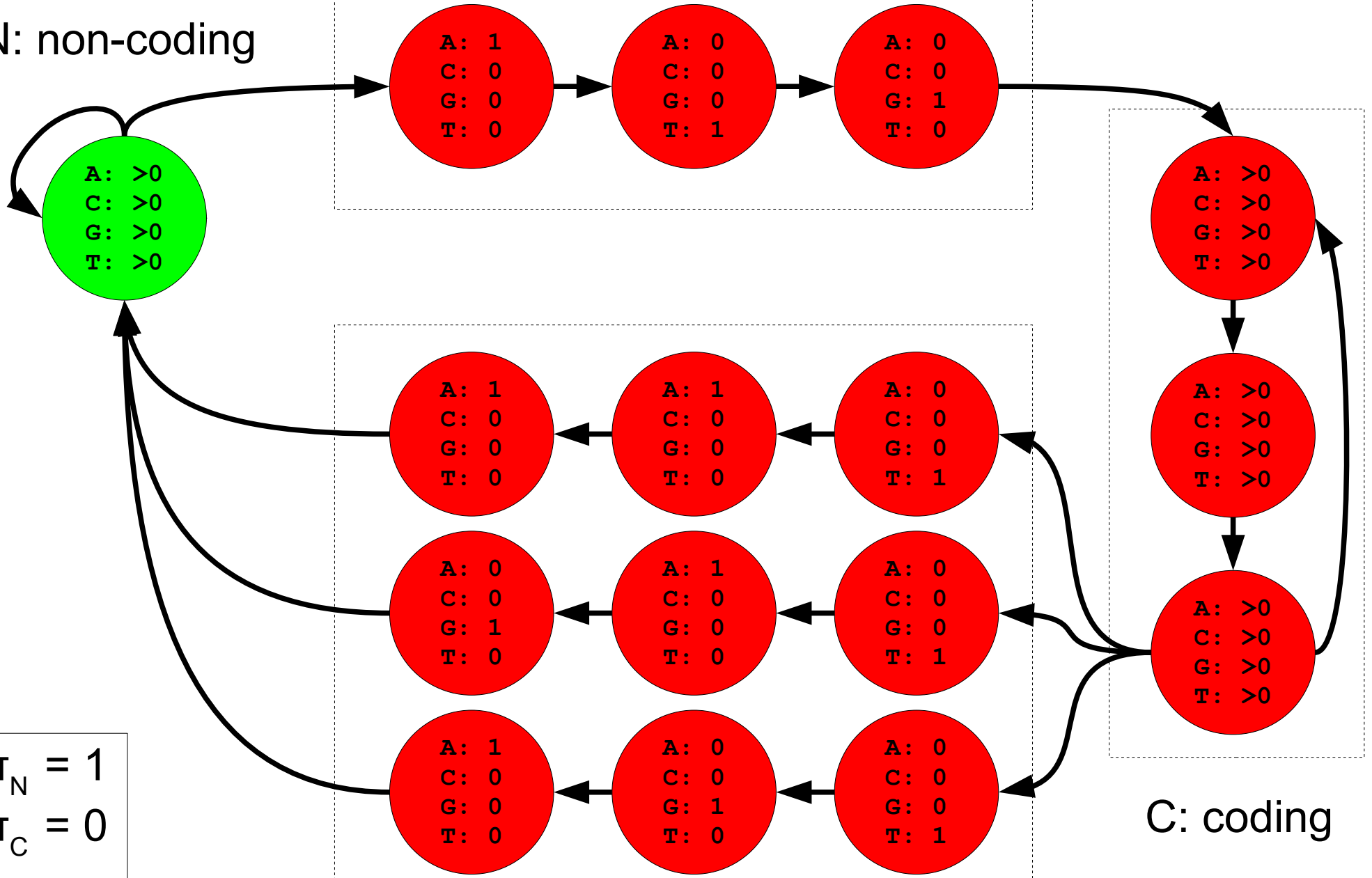
C: coding



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 multiple of 3

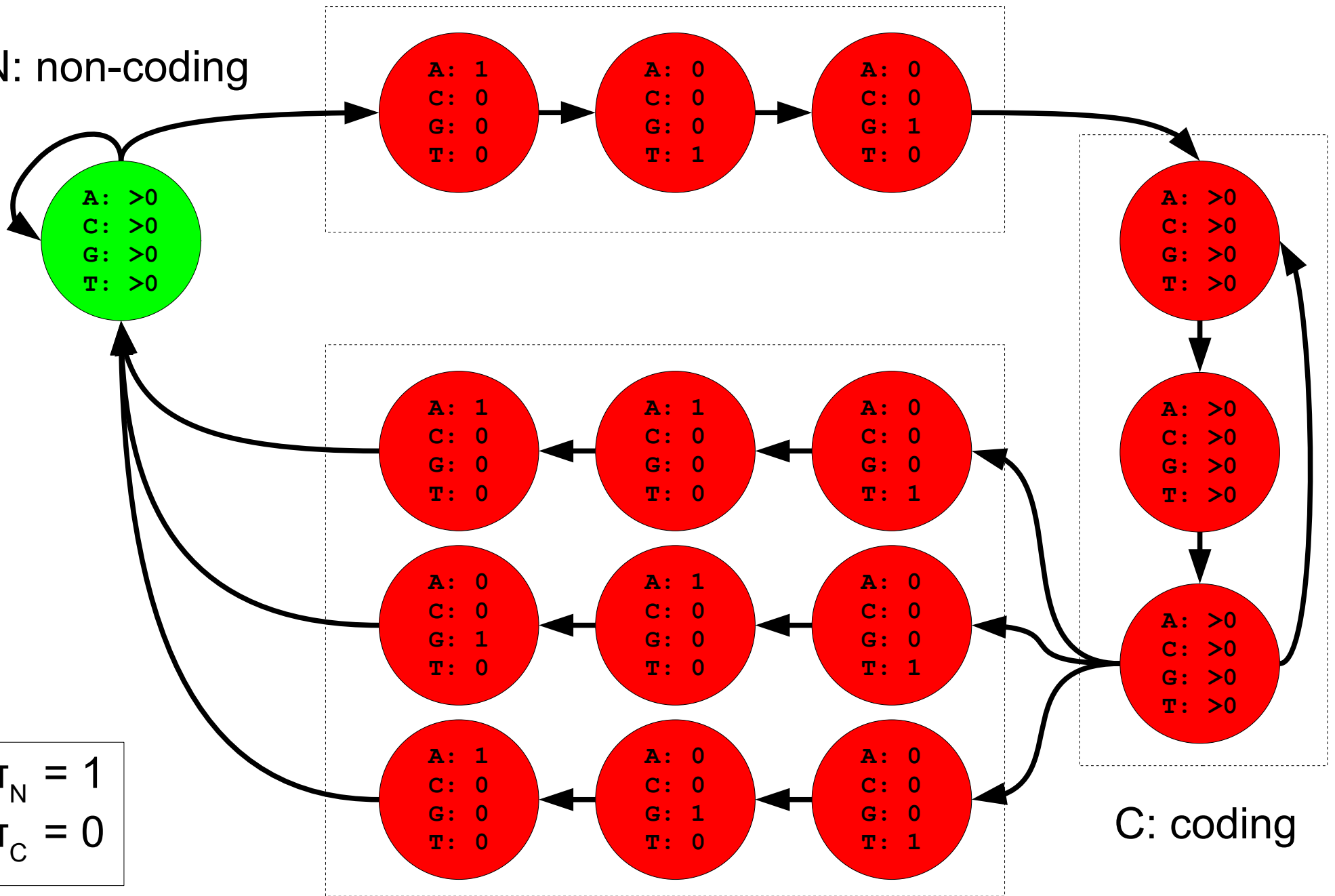
N: non-coding



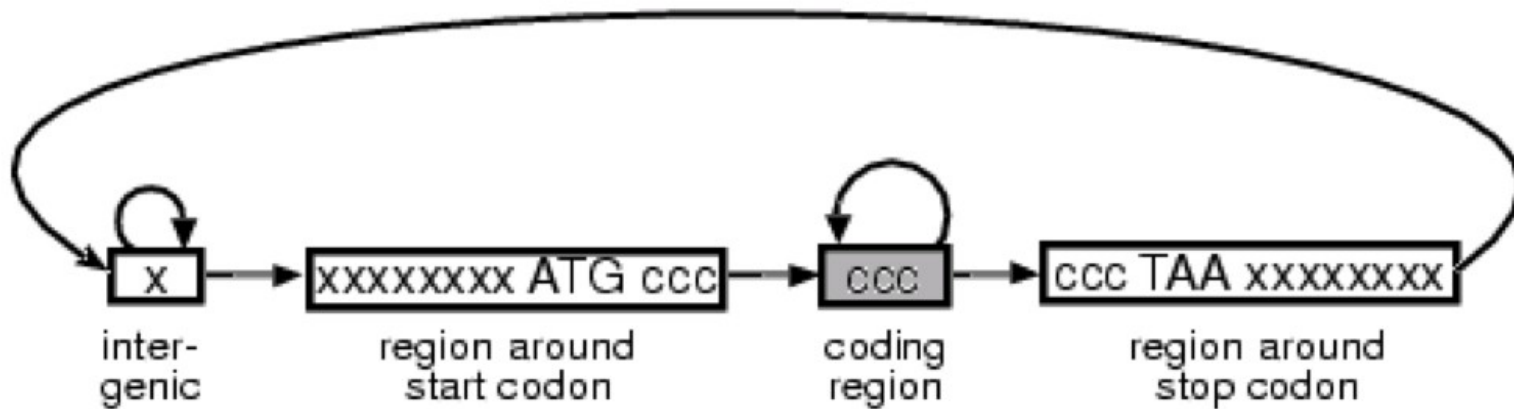
C: coding

Gene structure in procaryotes

N: non-coding



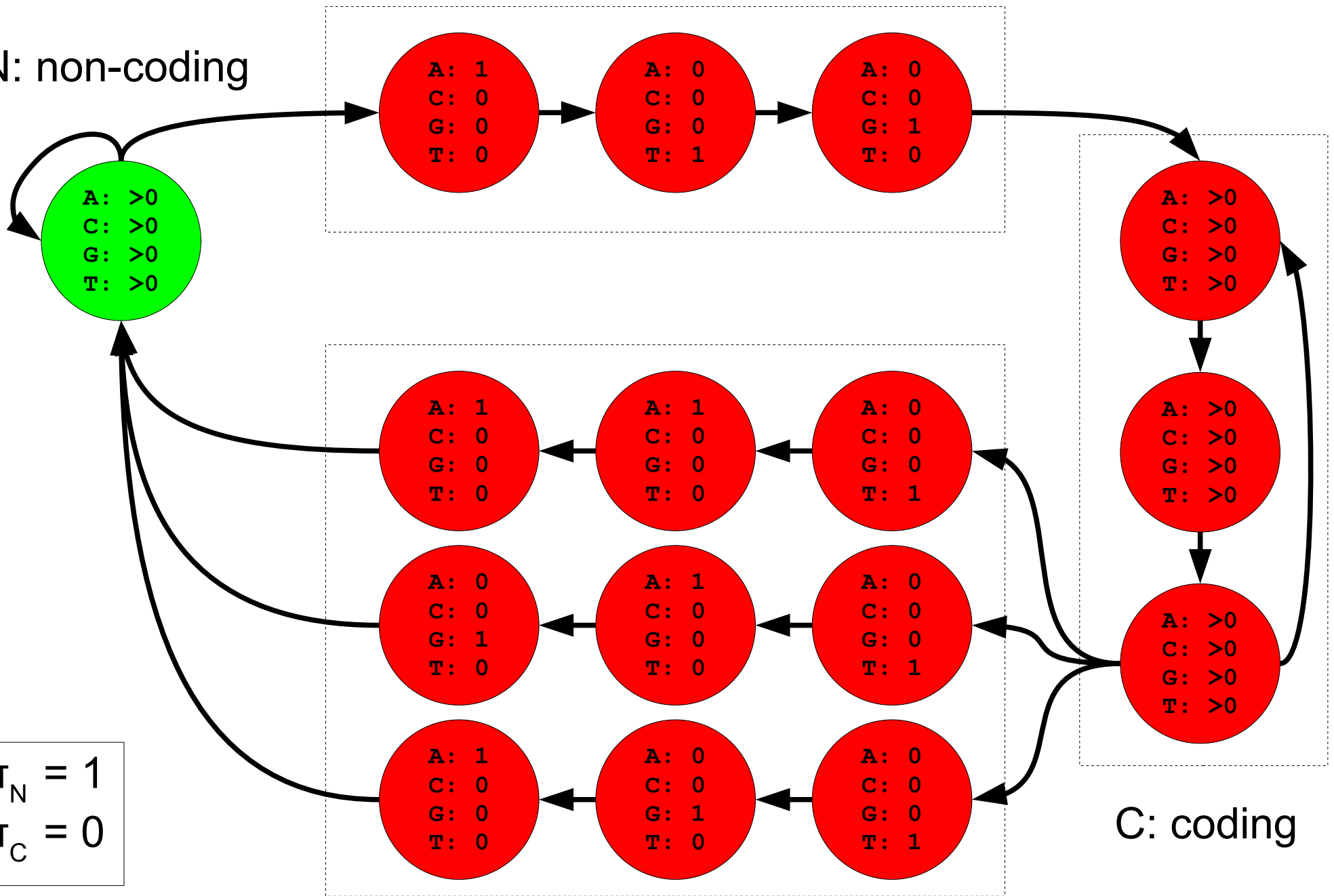
Gene structure in procaryotes



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

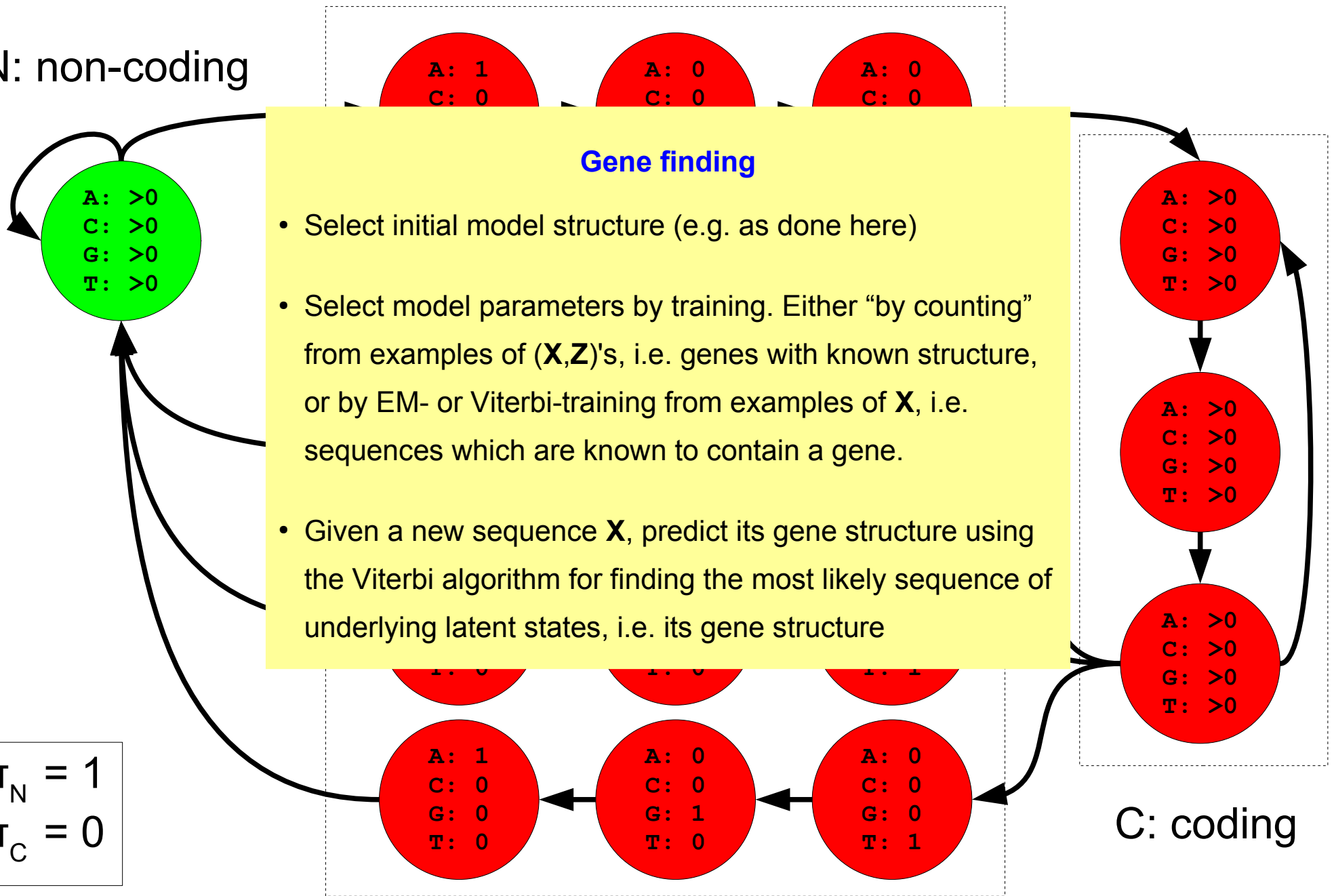
Gene structure in procaryotes

N: non-coding



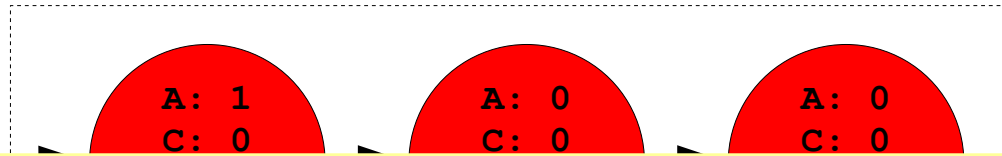
Gene structure in procaryotes

N: non-coding



Example – Gene finding


N: non-coding

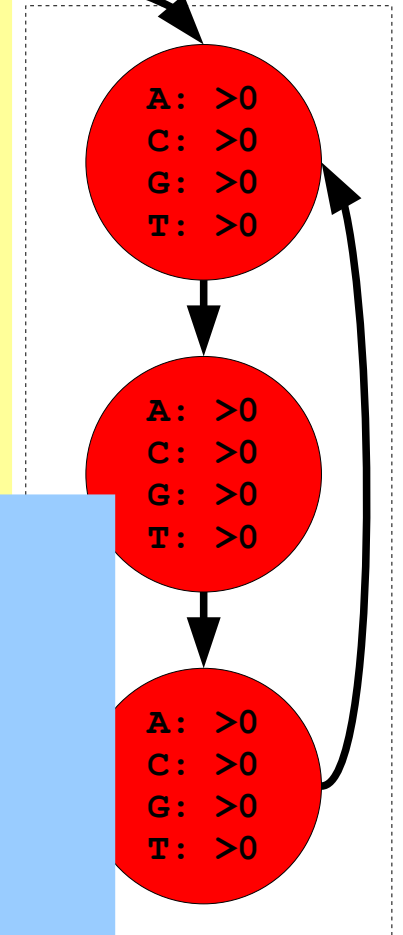


Gene finding

- Select initial model structure (e.g. as done here)
- Select model parameters by training. Either “by counting” from examples of (\mathbf{X}, \mathbf{Z}) 's, i.e. genes with known structure, or by EM- or Viterbi-training from examples of \mathbf{X} , i.e. sequences which are known to contain a gene.

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 ...

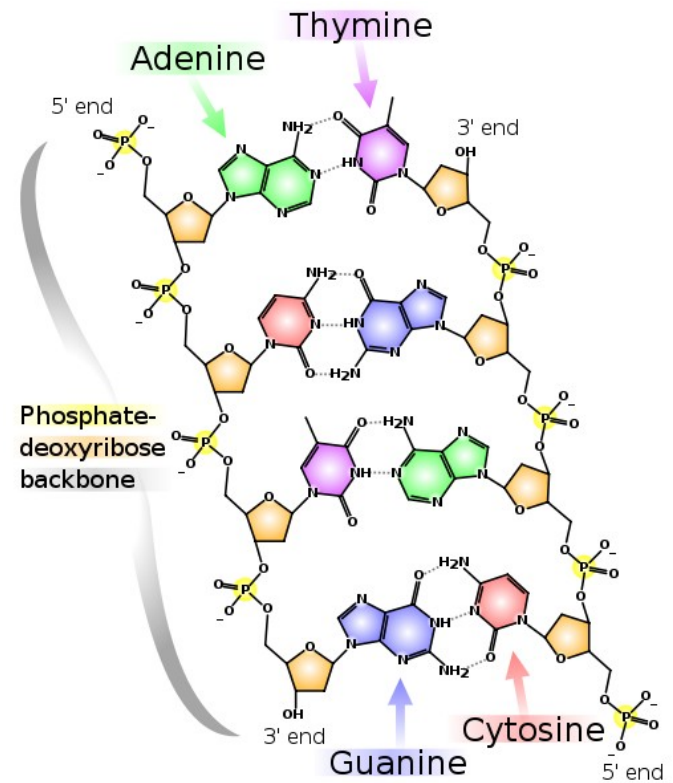
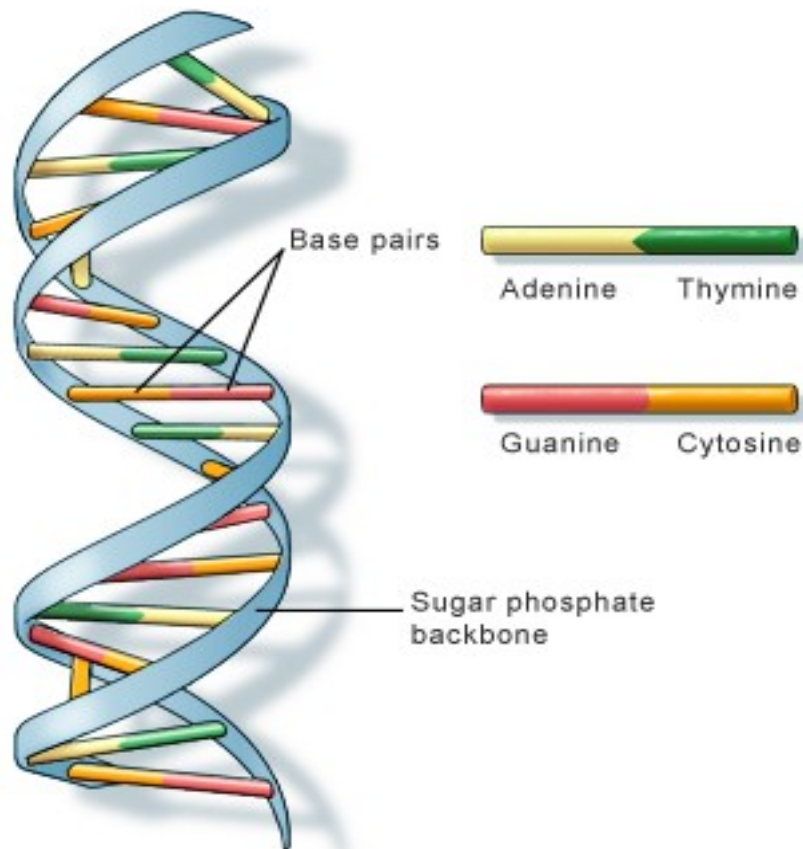
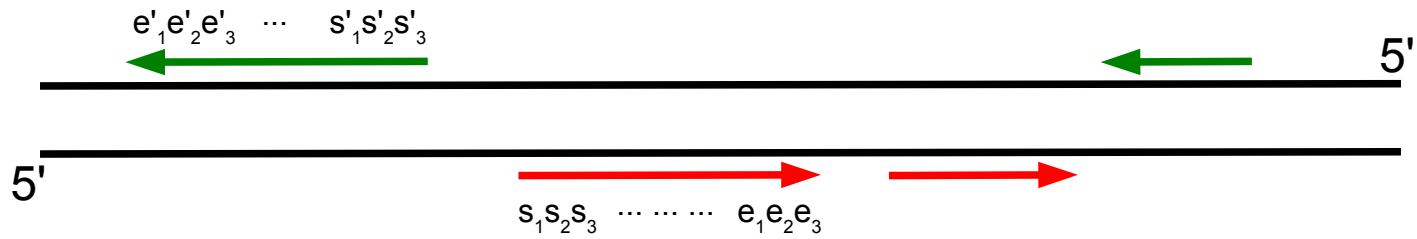


C: coding

$$\pi_N = 1$$

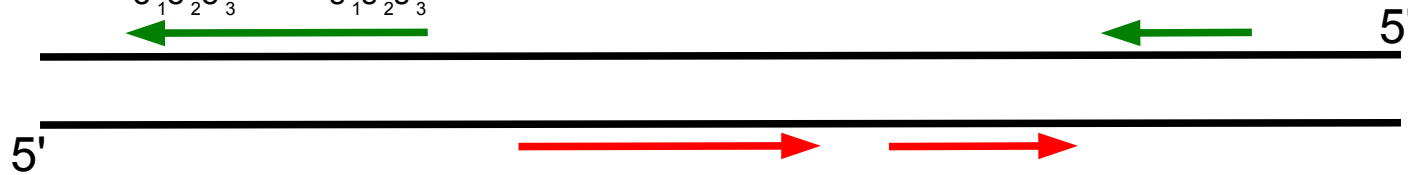
$$\pi_C = 0$$

DNA



DNA

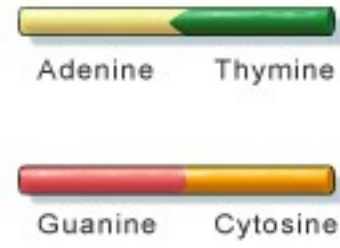
AGT
GAT
AAT GTA
e'₁e'₂e'₃ ... s'₁s'₂s'₃



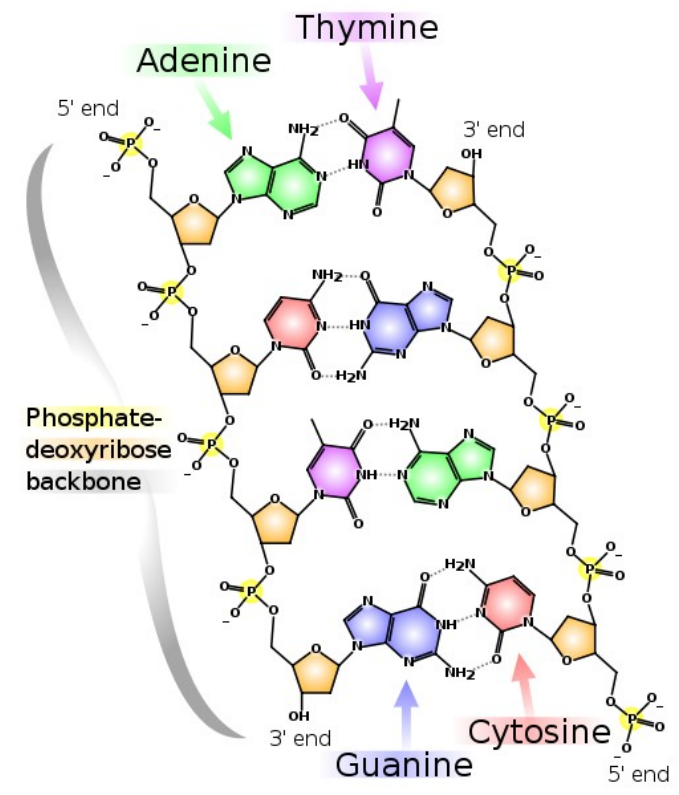
TTA CAT s₁s₂s₃ ... e₁e₂e₃
CTA **ATG** **TAA**
TCA **TAG**
 TGA



Base pairs



Sugar phosphate backbone



Even more biology

There can be genes in both directions



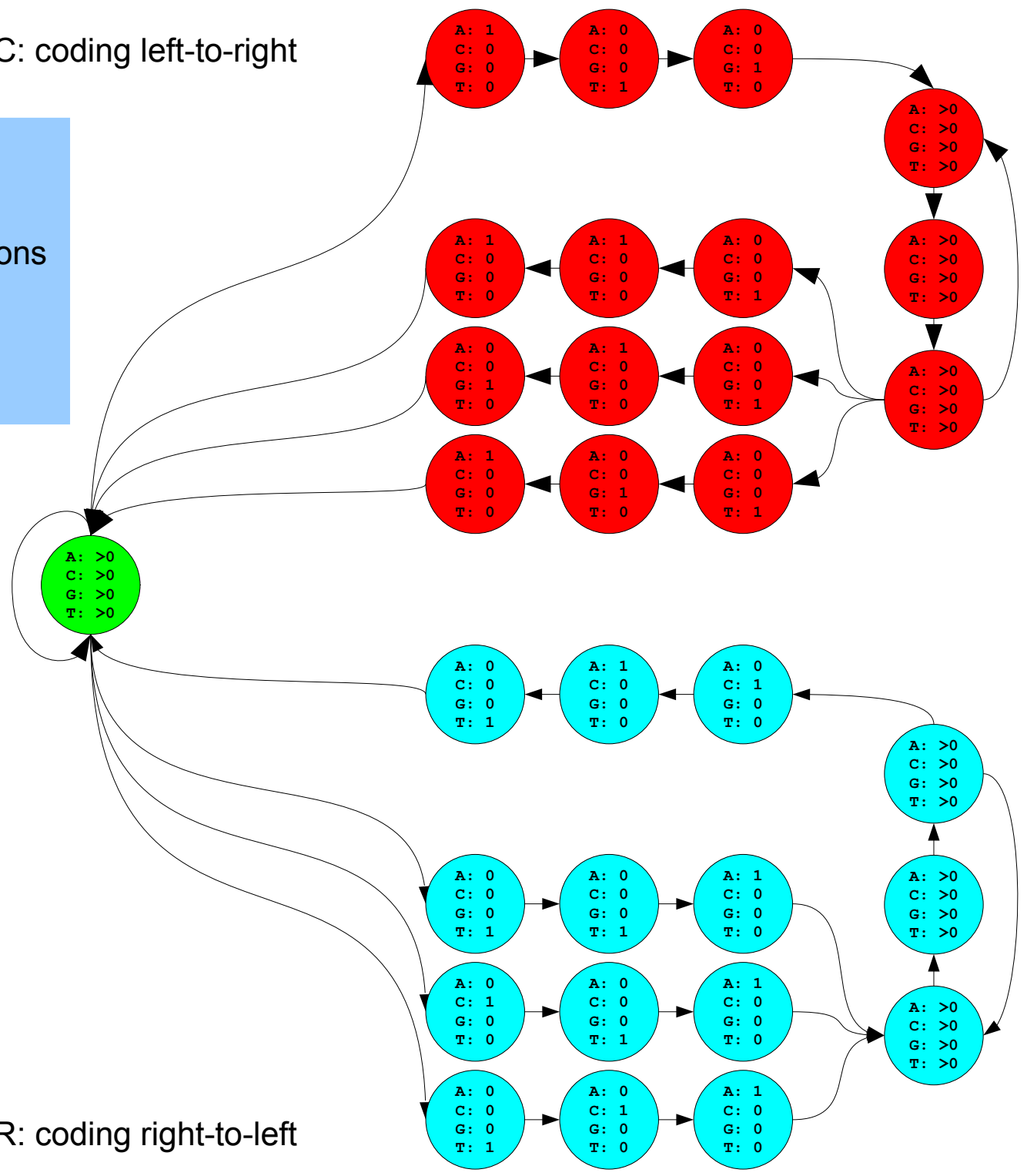
C: coding left-to-right

N: Non-coding

R: coding right-to-left

$$\pi_N = 1$$

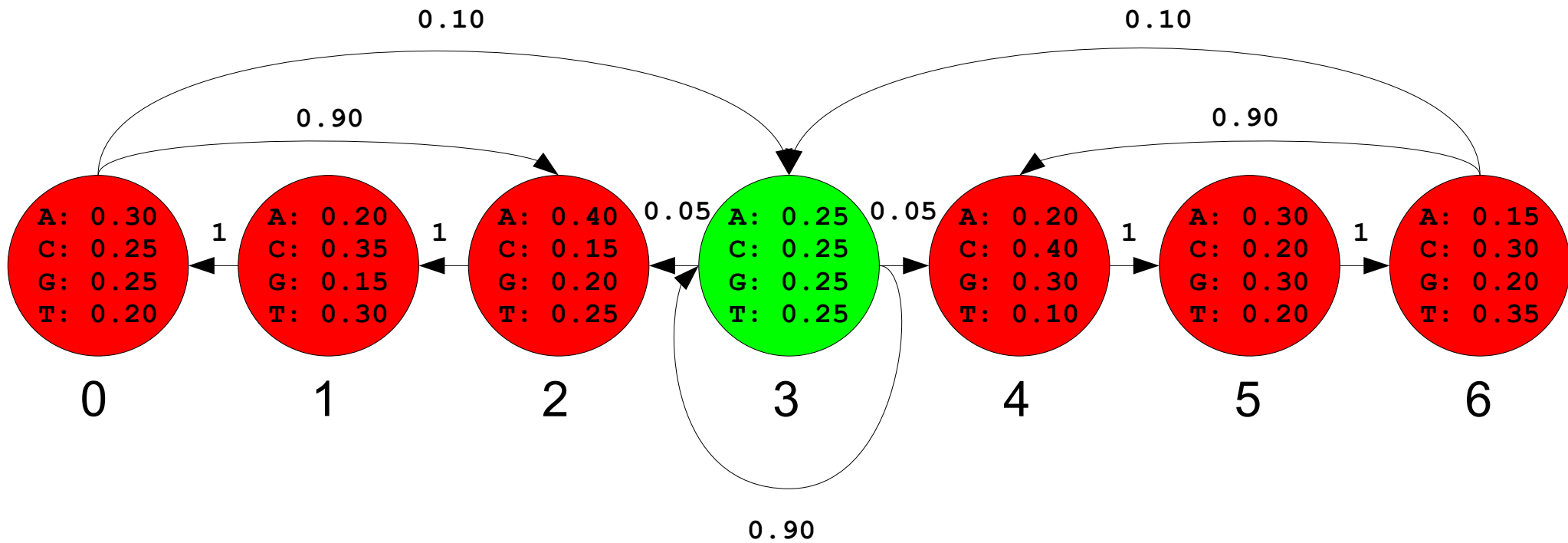
$$\pi_C = 0$$



Example – 7-state HMM

Observable: {A, C, G, T}, States: {0, 1, 2, 3, 4, 5, 6}

A	<table border="1" style="font-family: monospace; font-size: 0.8em;"> <tr><td>0.00</td><td>0.00</td><td>0.90</td><td>0.10</td><td>0.00</td><td>0.00</td><td>0.00</td></tr> <tr><td>1.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td></tr> <tr><td>0.00</td><td>1.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td></tr> <tr><td>0.00</td><td>0.00</td><td>0.05</td><td>0.90</td><td>0.05</td><td>0.00</td><td>0.00</td></tr> <tr><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>1.00</td><td>0.00</td></tr> <tr><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>0.00</td><td>1.00</td></tr> <tr><td>0.00</td><td>0.00</td><td>0.00</td><td>0.10</td><td>0.90</td><td>0.00</td><td>0.00</td></tr> </table>	0.00	0.00	0.90	0.10	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.05	0.90	0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.10	0.90	0.00	0.00	π	<table border="1" style="font-family: monospace; font-size: 0.8em;"> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> <tr><td>1.00</td></tr> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> <tr><td>0.00</td></tr> </table>	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	φ	<table border="1" style="font-family: monospace; font-size: 0.8em;"> <tr><td>0.30</td><td>0.25</td><td>0.25</td><td>0.20</td></tr> <tr><td>0.20</td><td>0.35</td><td>0.15</td><td>0.30</td></tr> <tr><td>0.40</td><td>0.15</td><td>0.20</td><td>0.25</td></tr> <tr><td>0.25</td><td>0.25</td><td>0.25</td><td>0.25</td></tr> <tr><td>0.20</td><td>0.40</td><td>0.30</td><td>0.10</td></tr> <tr><td>0.30</td><td>0.20</td><td>0.30</td><td>0.20</td></tr> <tr><td>0.15</td><td>0.30</td><td>0.20</td><td>0.35</td></tr> </table>	0.30	0.25	0.25	0.20	0.20	0.35	0.15	0.30	0.40	0.15	0.20	0.25	0.25	0.25	0.25	0.25	0.20	0.40	0.30	0.10	0.30	0.20	0.30	0.20	0.15	0.30	0.20	0.35
0.00	0.00	0.90	0.10	0.00	0.00	0.00																																																																																					
1.00	0.00	0.00	0.00	0.00	0.00	0.00																																																																																					
0.00	1.00	0.00	0.00	0.00	0.00	0.00																																																																																					
0.00	0.00	0.05	0.90	0.05	0.00	0.00																																																																																					
0.00	0.00	0.00	0.00	0.00	1.00	0.00																																																																																					
0.00	0.00	0.00	0.00	0.00	0.00	1.00																																																																																					
0.00	0.00	0.00	0.10	0.90	0.00	0.00																																																																																					
0.00																																																																																											
0.00																																																																																											
0.00																																																																																											
1.00																																																																																											
0.00																																																																																											
0.00																																																																																											
0.00																																																																																											
0.00																																																																																											
0.00																																																																																											
0.30	0.25	0.25	0.20																																																																																								
0.20	0.35	0.15	0.30																																																																																								
0.40	0.15	0.20	0.25																																																																																								
0.25	0.25	0.25	0.25																																																																																								
0.20	0.40	0.30	0.10																																																																																								
0.30	0.20	0.30	0.20																																																																																								
0.15	0.30	0.20	0.35																																																																																								



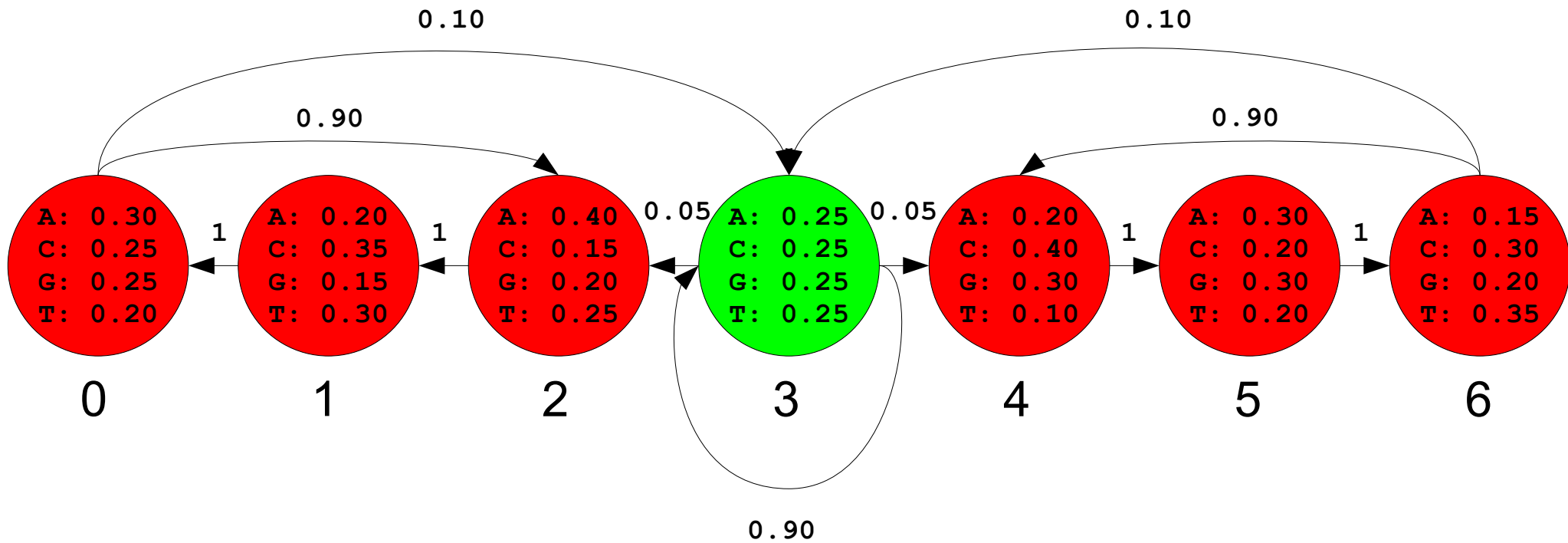
This model is also applicable for gene finding.

It does not model start- and stop-codons explicitly, but models that genes in both directions are a sequence of triplets.

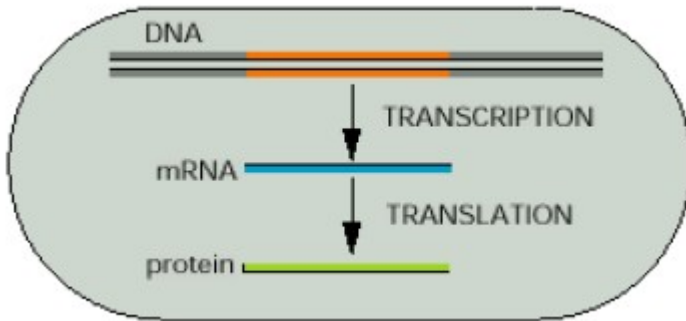
<i>A</i>	0.00	0.00	0.90	0.10	0.00	0.00	0.00
	1.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	1.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.05	0.90	0.05	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	1.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	1.00
	0.00	0.00	0.00	0.10	0.90	0.00	0.00
	0.00	0.00	0.00	0.10	0.90	0.00	0.00

π	0.00
	0.00
	0.00
	1.00
	0.00
	0.00
	0.00
	0.00

φ	0.30	0.25	0.25	0.20
	0.20	0.35	0.15	0.30
	0.40	0.15	0.20	0.25
	0.25	0.25	0.25	0.25
	0.20	0.40	0.30	0.10
	0.30	0.20	0.30	0.20
	0.15	0.30	0.20	0.35
	0.15	0.30	0.20	0.35



Problem: From annotation to Z



Biological facts

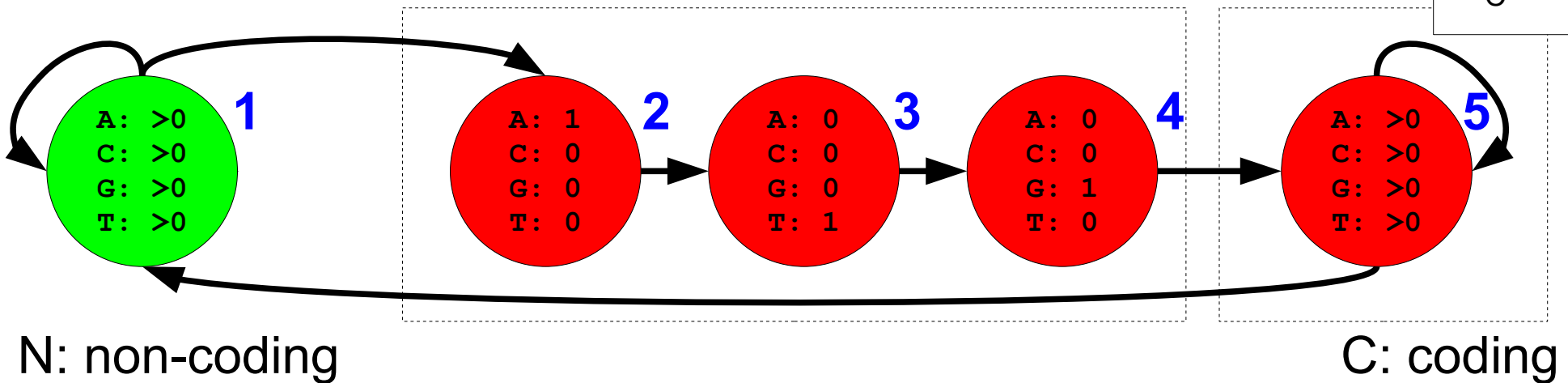
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$



Problem: From annotation to Z

Problem: The string $Z=NNNCCC\dots$ is not a proper sequence of states in the illustrated HMM, but it can easily be converted into one (because there in this case is a 1-1 matching between a sequence of Ns and Cs and a sequence of states).

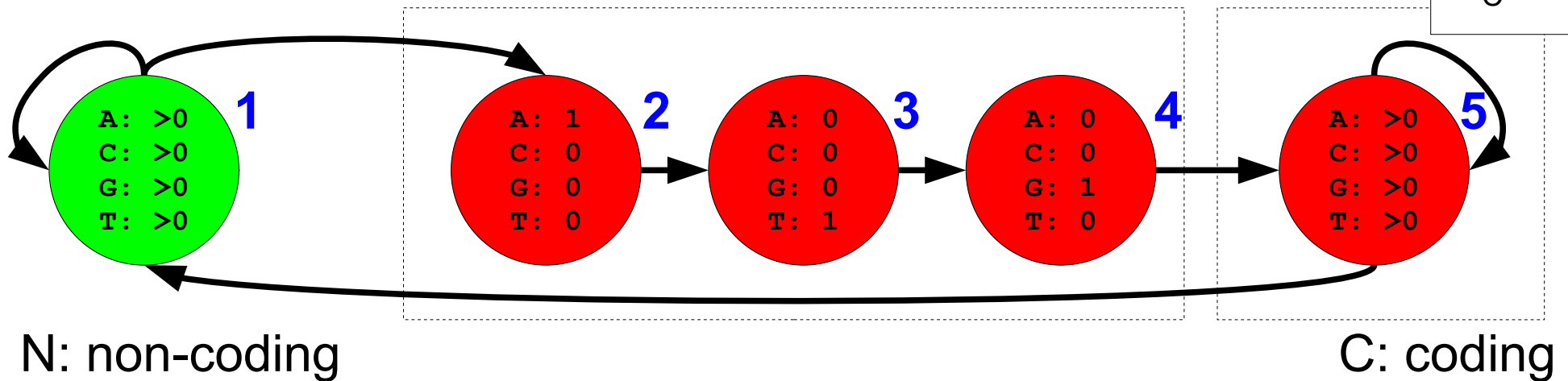
ence of A,C,G,T's

Z: NNNCCCCCCCCNNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$



Problem: From annotation to Z

Problem: The string $Z=NNNCCC\dots$ is not a proper sequence of states in the illustrated HMM, but it can easily be converted into one (because there is a 1-1 matching between a sequence of Ns and Cs and a sequence of states).

ence of A,C,G,T's

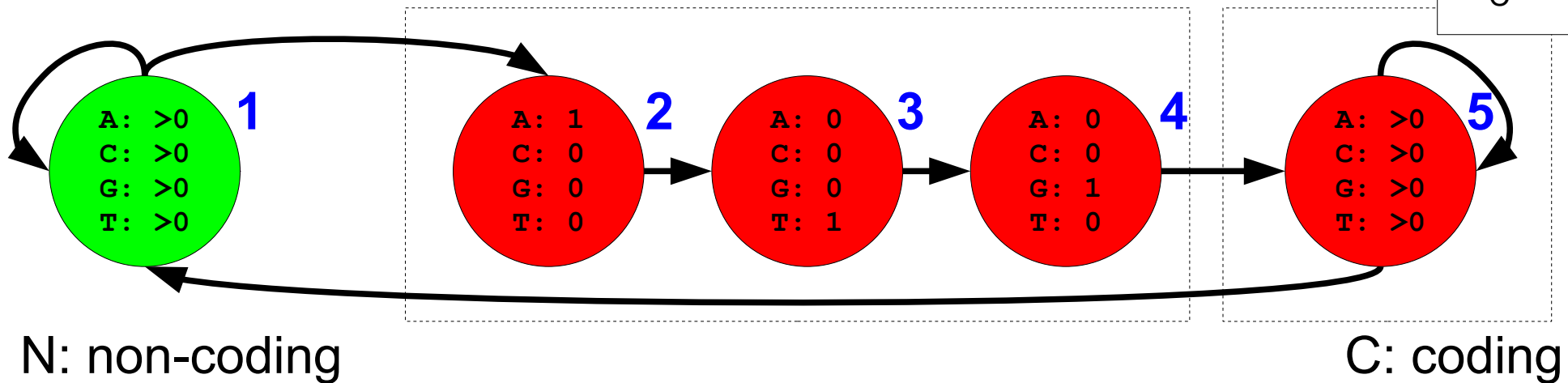
1112345555511111111234555555555555511111111111

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

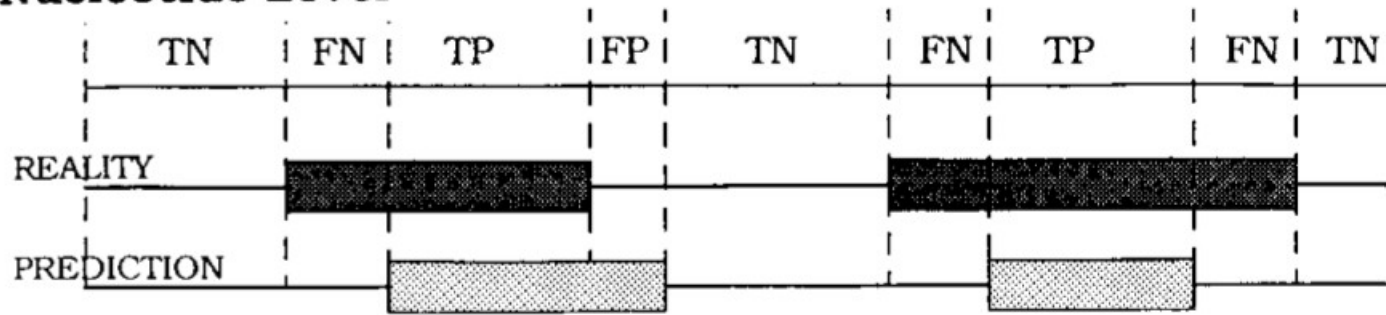
$$\pi_N = 1$$

$$\pi_C = 0$$



Evaluating performance

Nucleotide Level



		REALITY		
		coding	no coding	
PREDICTION	coding	TP	FP	TP+FP
	no coding	FN	TN	FN+TN
		TP+FN	TN+FP	

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

Sensitivity

$$S_p = \frac{TN}{TN + FP}$$

Specificity

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

Correlation Coefficient

$$ACP = \frac{1}{4} \left[\frac{TP}{TP + FN} + \frac{TP}{TP + FP} + \frac{TN}{TN + FP} + \frac{TN}{TN + FN} \right]$$

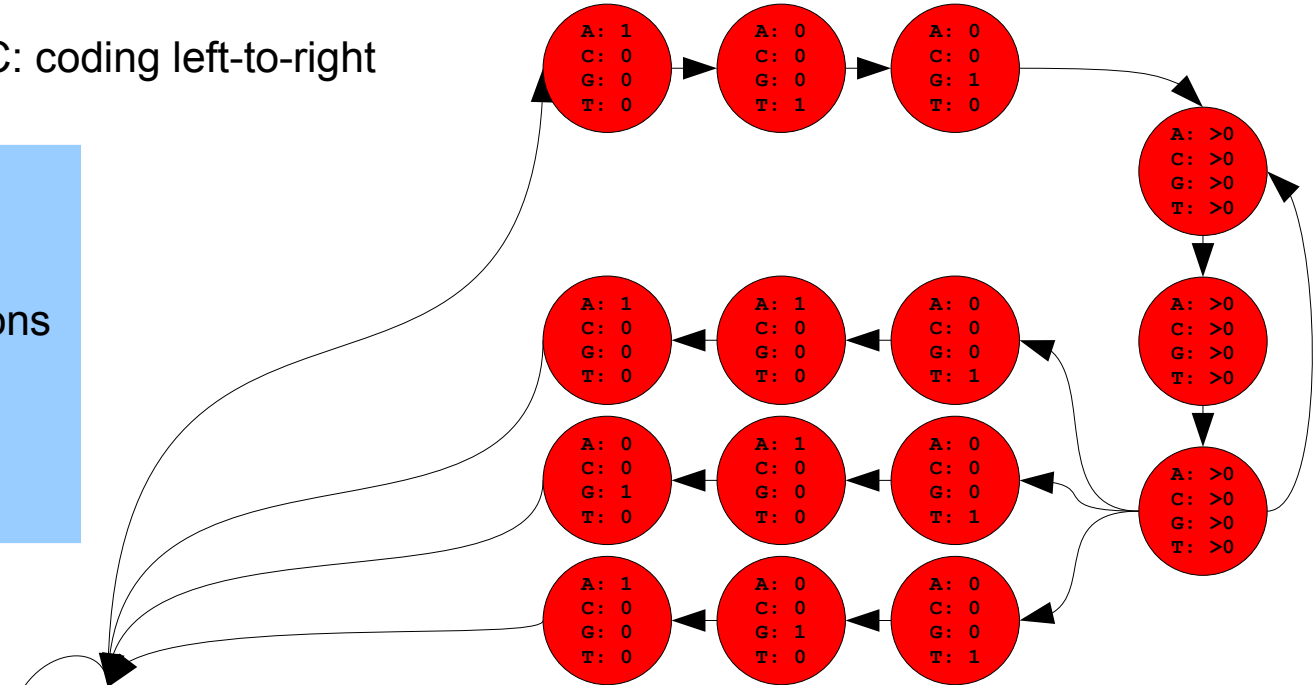
$$AC = (ACP - 0.5) \times 2$$

Approximate Correlation

Even more biology

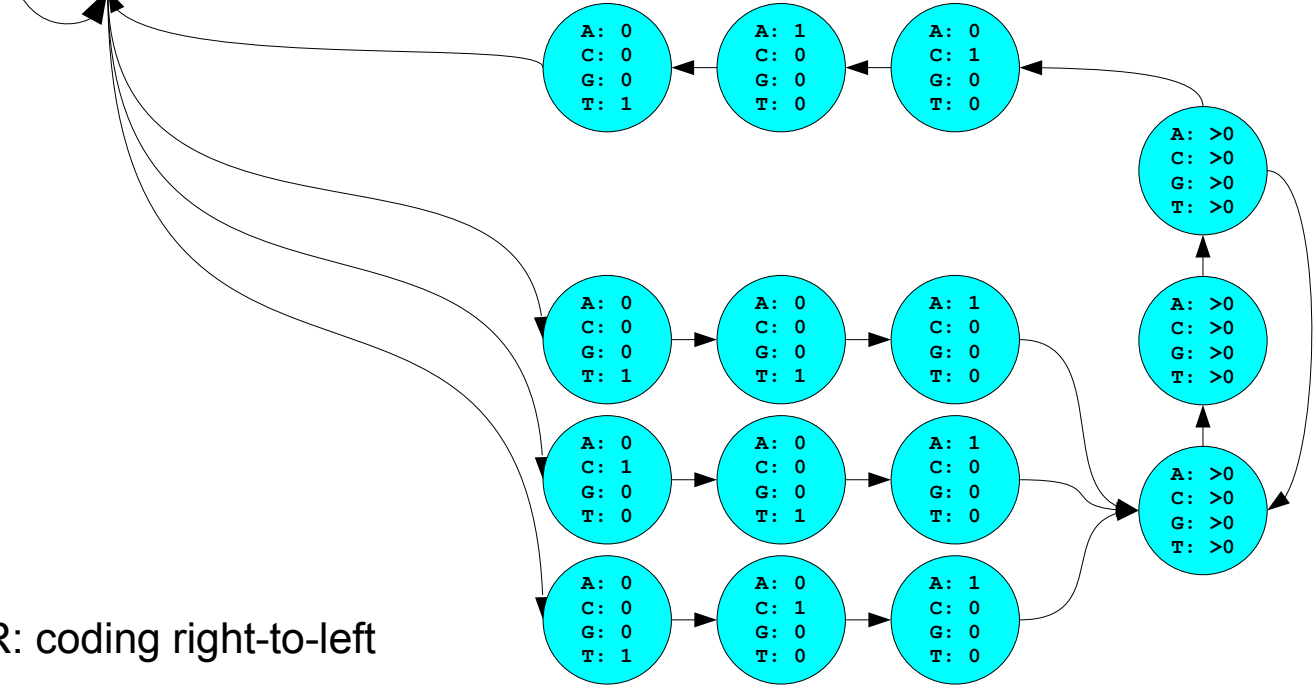
There can be genes in both directions

C: coding left-to-right



N: Non-coding

R: coding right-to-left



$$\pi_N = 1$$

$$\pi_C = 0$$

Analysis of some genomes

```
Length of genome1: 1852441 (1852441)
Length of genome2: 2211485 (2211485)
Length of genome3: 2499279 (2499279)
Length of genome4: 1796846 (1796846)
Length of genome5: 2685015 (2685015)
Length of genome6: 2127839 (2127839)
Length of genome7: 2742531 (2742531)
Length of genome8: 2046115 (2046115)
Length of genome9: 2388435 (2388435)
Length of genome10: 1570485 (1570485)
Length of genome11: 2096309 (2096309)
```

Start-codon in normal genes:

```
ATG [8423, 'NCCC']
ATC [3, 'NCCC']
ATA [1, 'RCCC']
GTG [713, 'NCCC']
ATT [3, 'NCCC']
CTG [2, 'NCCC']
GTT [1, 'NCCC']
CTC [1, 'NCCC']
TTA [1, 'NCCC']
TTG [1020, 'NCCC']
```

Stop-codon in normal genes:

```
TAG [1949, 'CCCN']
TGA [1531, 'CCCN']
TAA [6686, 'CCCN']
```

Reversed stop-codon in reversed genes:

```
TTA (reverse-complement: TAA) [6596, 'NRRR']
CTA (reverse-complement: TAG) [2014, 'NRRR']
TCA (reverse-complement: TGA) [1148, 'NRRR']
```

Reversed start-codon in reversed genes:

```
TAT (reverse-complement: ATA) [2, 'RRRN']
ATG (reverse-complement: CAT) [1, 'RRRN']
GAT (reverse-complement: ATC) [1, 'RRRN']
CAT (reverse-complement: ATG) [8077, 'RRRN']
AAT (reverse-complement: ATT) [4, 'RRRN']
TAC (reverse-complement: GTA) [1, 'RRRN']
CAC (reverse-complement: GTG) [715, 'RRRN']
CAA (reverse-complement: TTG) [953, 'RRRN']
CAG (reverse-complement: CTG) [4, 'RRRN']
```