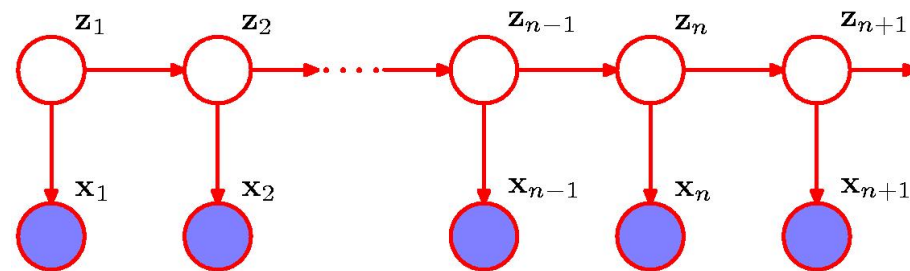


Hidden Markov Models

Selecting the initial model parameters

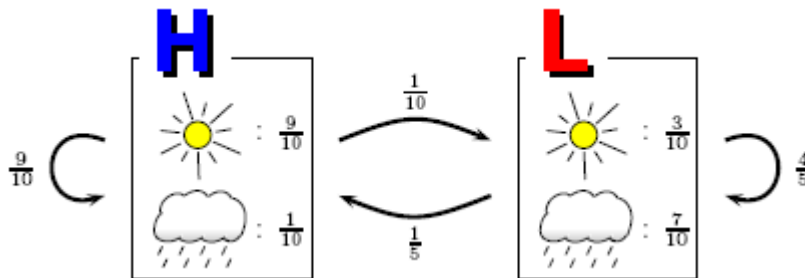
Using HMMs for (simple) gene finding



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

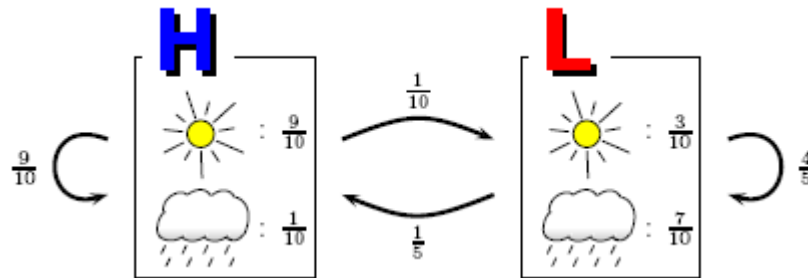
And *emits* a sequence of symbols:



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 ...

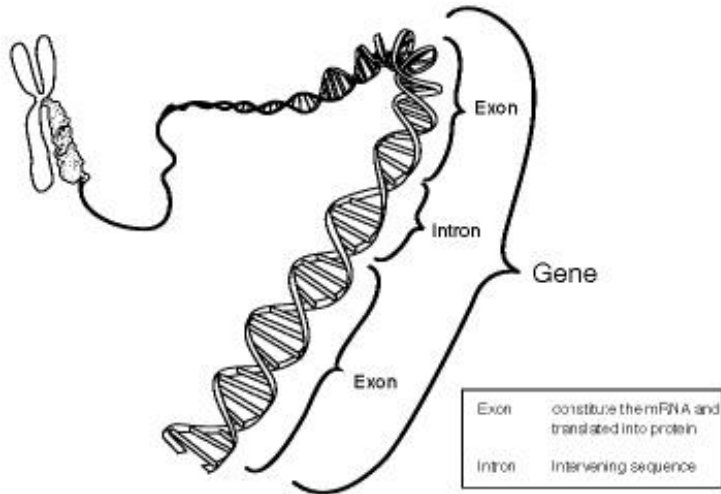
Selecting initial model parameters

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.



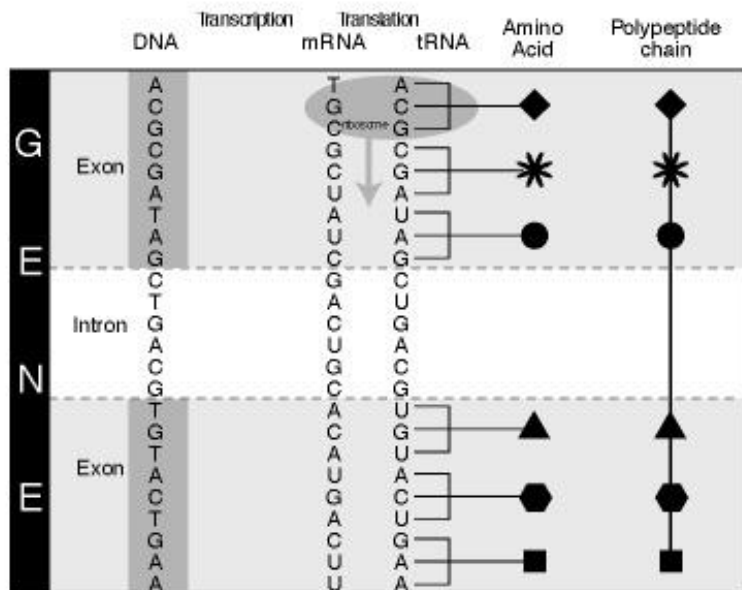
The initial selection of parameters is essential just to decide which parameters are 0 (or 1), i.e. to decide which transitions or emissions should never (or always) be possible ...

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

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

>NC_002737.1 Streptococcus pyogenes M1 GAS

TTGTTGATATTCTGTTTTTCTTTTTTAGTTTTCCACATGAAAAATAGTTGAAAACAATA
GCGGTGTCCCCTTAAAATGGCTTTTCCACAGGTTGTGGAGAACCCAAATTAACAGTGTTA
ATTTATTTTCCACAGGTTGTGGAAAACTAACTATTATCCATCGTTCTGTGGAAAACTAG
AATAGTTTATGGTAGAATAGTTCTAGAATTATCCACAAGAAGGAACCTAGTATGACTGAA
AATGAACAAATTTTTTGGAACAGGGTCTTGGAATTAGCTCAGAGTCAATTAACAGGCA
ACTTATGAATTTTTTGTTCATGATGCCCCTCTATTAAGGTCGATAAGCATATTGCAACT
ATTTACTTAGATCAAATGAAAGAGCTCTTTTGGGAAAAAATCTTAAAGATGTTATTCTT
ACTGCTGGTTTTGAAGTTTATAACGCTCAAATTTCTGTTGACTATGTTTTCGAAGAAGAC
CTAATGATTGAGCAAAATCAGACCAAAATCAACCAAAAACCTAAGCAGCAAGCCTTAAAT
TCTTTGCCTACTGTTACTTCAGATTTAACTCGAAATATAGTTTTGAAAACCTTTATTCAA
GGAGATGAAAATCGTTGGGCTGTTGCTGCTTCAATAGCAGTAGCTAATACTCCTGGAAC
ACCTATAATCCTTTGTTTATTTGGGGTGGCCCTGGGCTTGGAACCCATTTATTAAT
GCTATTGGTAATTCTGTACTATTAGAAAATCCAAATGCTCGAATTAAATATATCACAGCT
GAAAACCTTTATTAATGAGTTTGTTATCCATATTCGCCTTGATACCATGGATGAATTGAAA
GAAAAATTTCTGAATTTAGATTTACTCCTTATTGATGATATCCAATCTTTAGCTAAAAAA
ACGCTCTCTGGAACACAAGAAGAGTTCTTTAATACTTTTAAATGCACTTCATAATAAAC
AAACAAATTGTCCTAACAAGCGACCGTACACCAGATCATCTCAATGATTTAGAAGATCGA
TTAGTTACTCGTTTTAAATGGGGATTAAACAGTCAATATCACACCTCCTGATTTTGAAACA
CGAGTGGCTATTTTGACAAATAAAATTCAAGAATATAACTTTATTTTTCCTCAAGATACC
ATTGAGTATTTGGCTGGTCAATTTGATTTCTAATGTCAGAGATTTAGAAGGTGCCTTAAAA
GATATTAGTCTGGTTGCTAATTTCAAACAAATTGACACGATTACTGTTGACATTGCTGCC
GAAGCTATTCGCGCCAGAAAGCAAGATGGACCTAAAATGACAGTTATTCCCATCGAAGAA
ATTCAAGCGCAAGTTGGAAAATTTACGGTGTTACCGTCAAAGAAATTAAGCTACTAAA
CGAACACAAAATATTGTTTTAGCAAGACAAGTAGCTATGTTTTTAGCACGTGAAATGACA
GATAACAGTCTTCCTAAAATTGGAAAAGAATTTGGTGGCAGAGACCATTCAACAGTACTC
CATGCCTATAATAAAATCAAAAACATGATCAGCCAGGACGAAAGCCTTAGGATCGAAATT
GAAACCATAAAAAACAAAATTAATAACATGTGGAAAAGAATATCTTTTATGAAATAGTT
ATCCACAAGTTGTGAACATCCATTTAGTCTTGGATTCTCTCGTTTATTTAGAGTTATCCA
CTATATACACAAGACCTACTACTACTATTATTATACTTATTAAATAAAGGAGTTCT

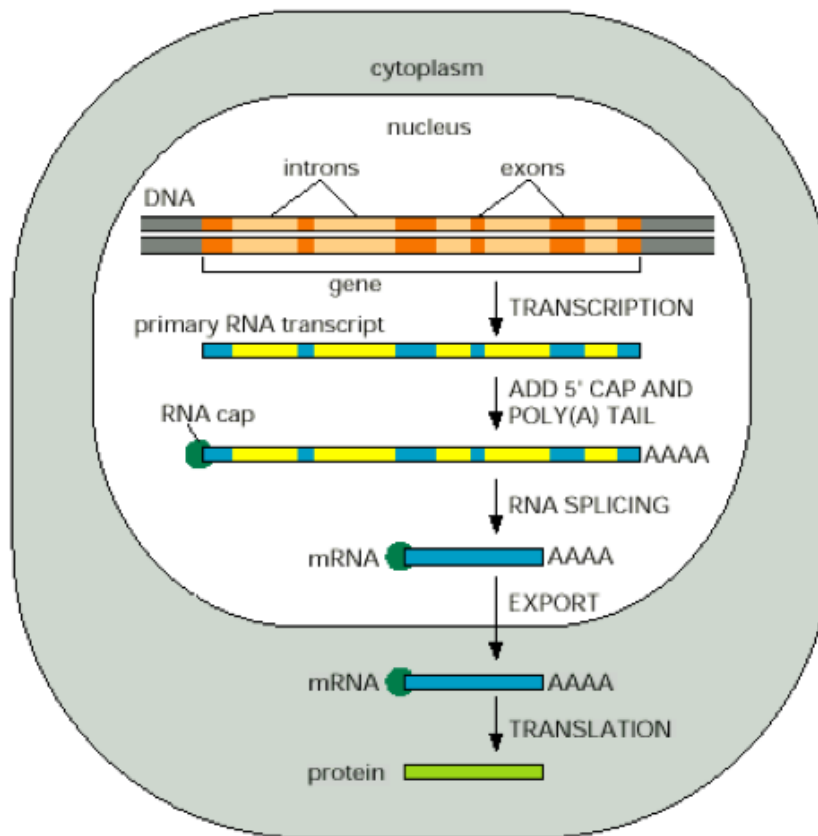
[illegible]

Design a HMM that models the syntax of genes

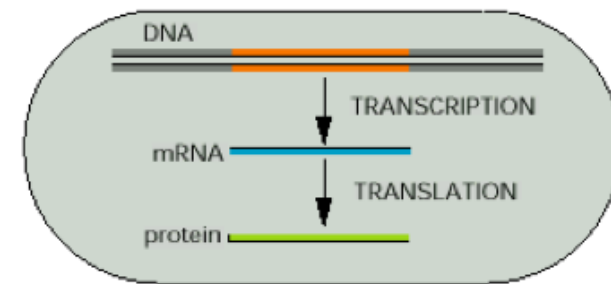
Gene structure

Depends on the organism (eucaryote or procaryote)

(A) EUCARYOTES



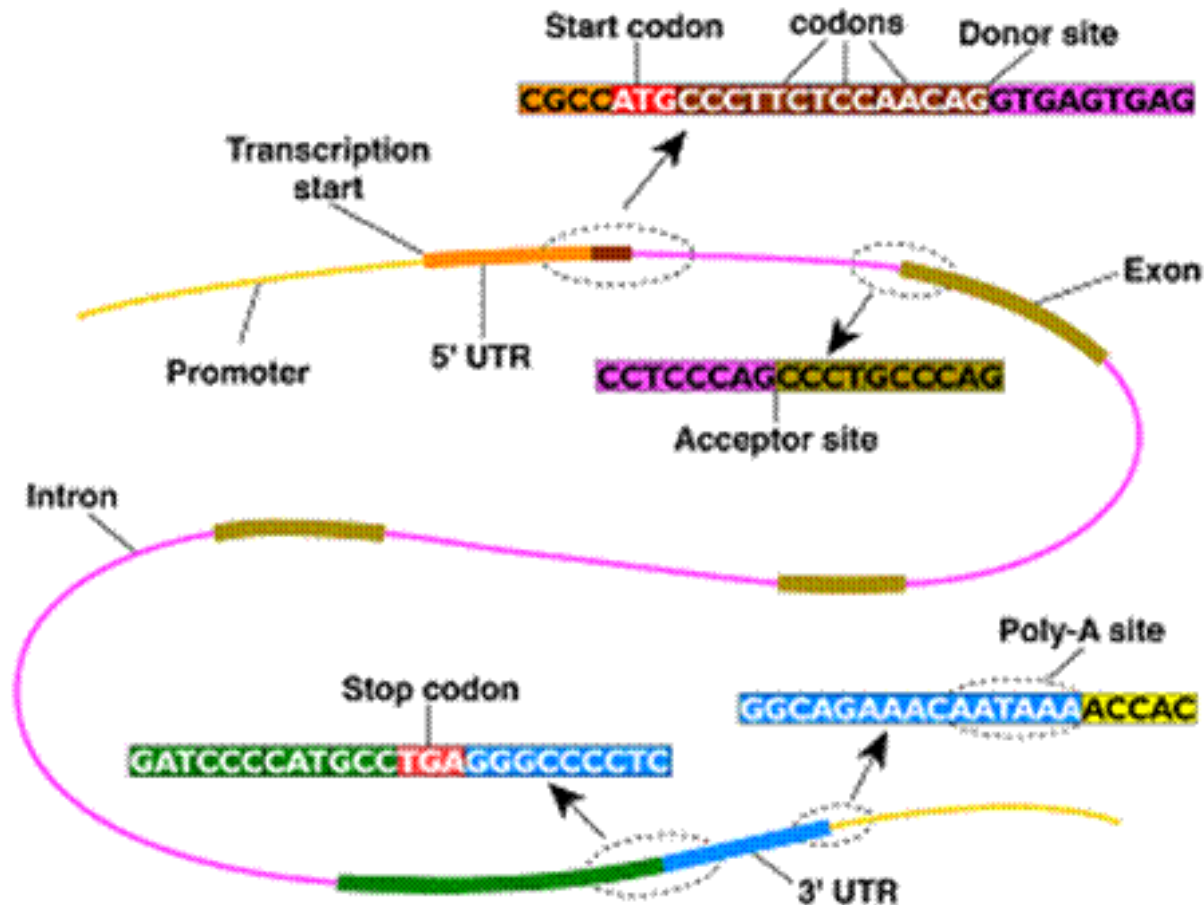
(B) PROCARYOTES



Smaller genomes and high coding density.

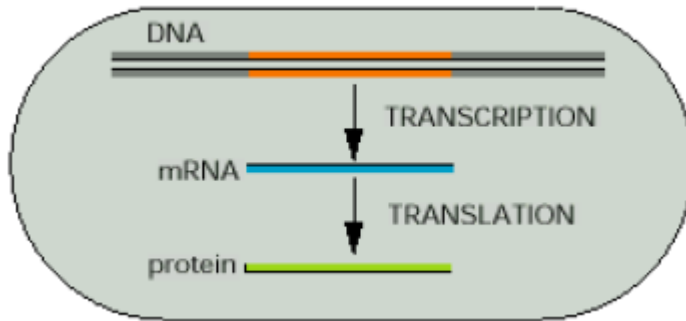
Large genomes. Intron/exon structure and low coding density

Gene structure in eukaryotes



Eukaryotic gene structure in more details

Gene structure in procaryotes

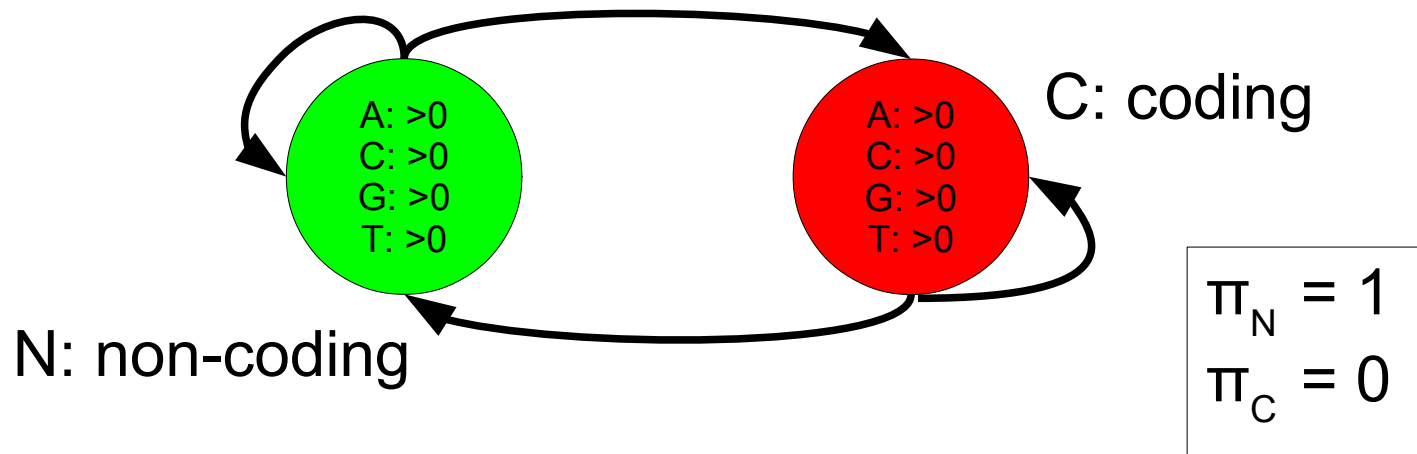


Biological facts

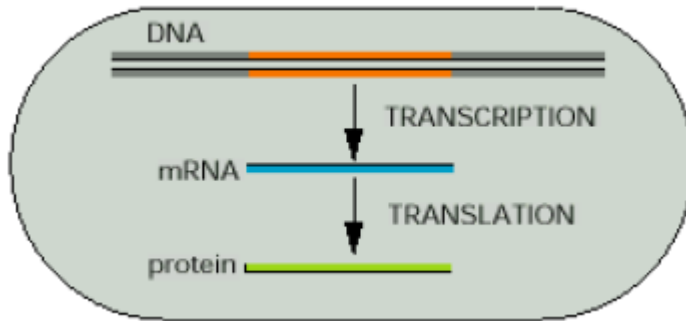
- The gene is a substring of the DNA sequence of A,C,G,T's

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacatgcag



Gene structure in procaryotes

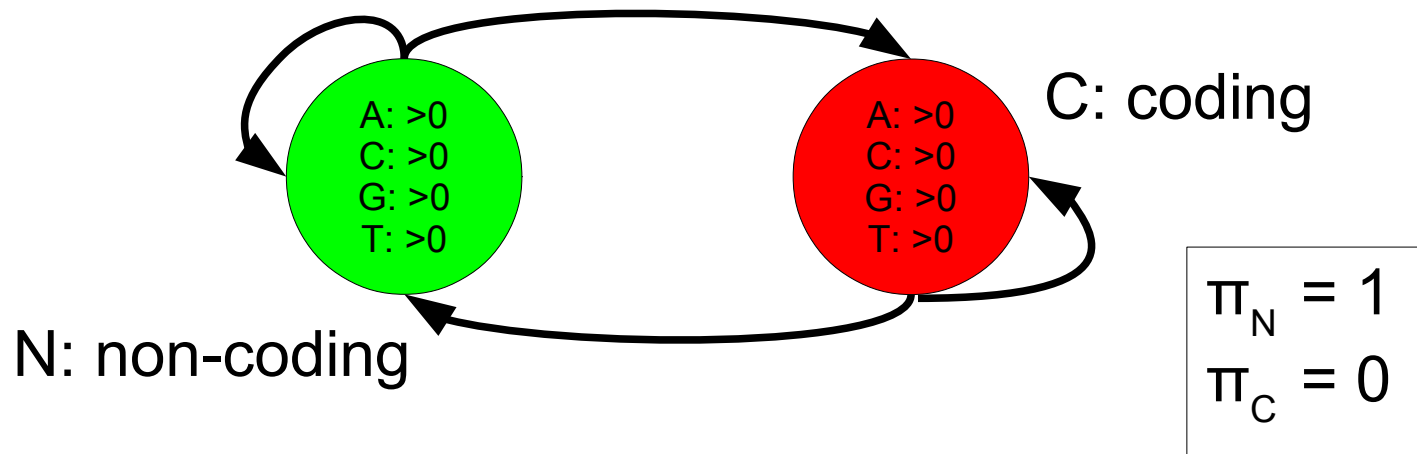


Biological facts

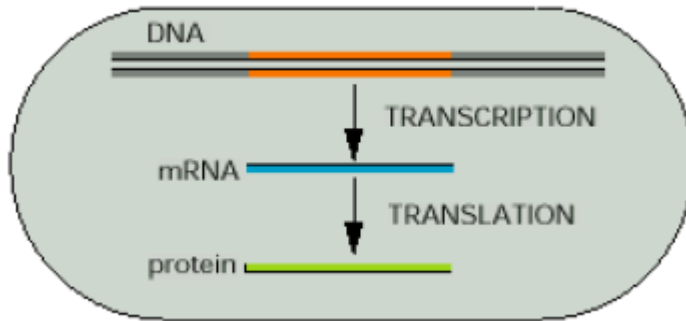
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacatgcag



Gene structure in procaryotes



Biological facts

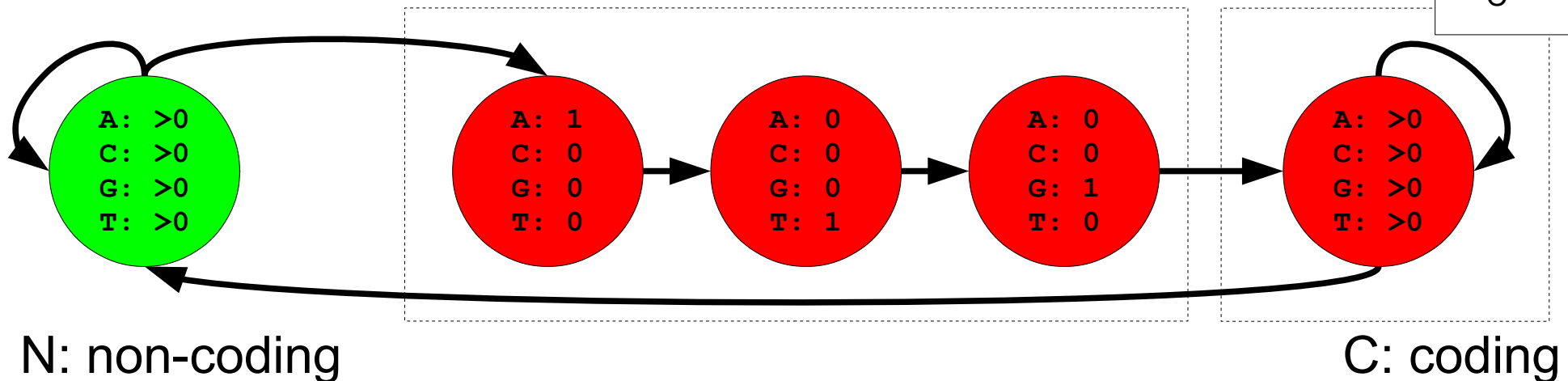
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

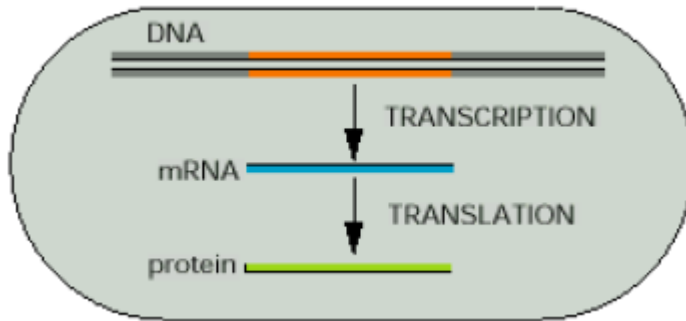
X: acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$



Gene structure in procaryotes



Biological facts

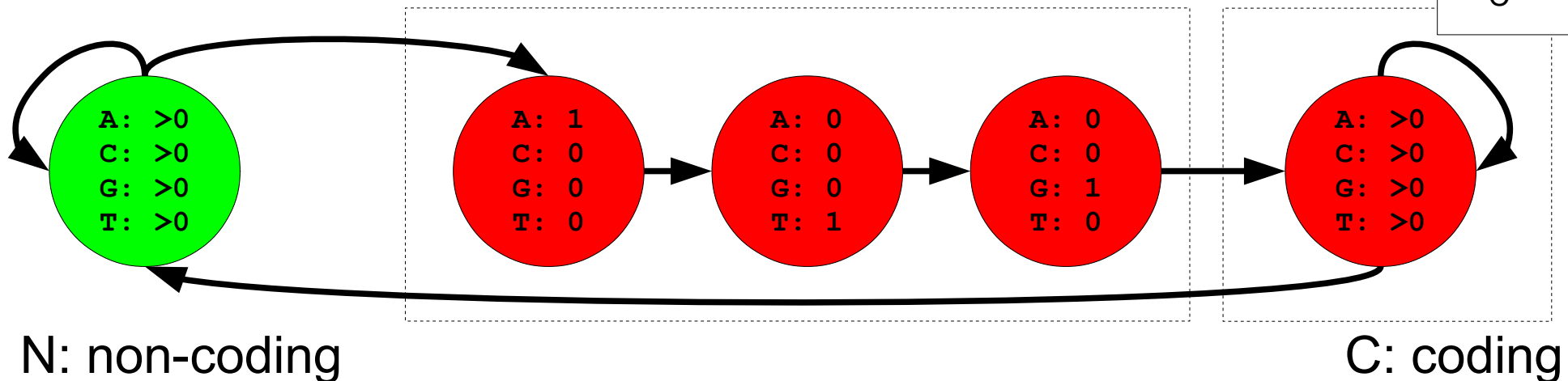
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**
- The gene ends with a stop-codon **taa**, **tag** or **tga**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$

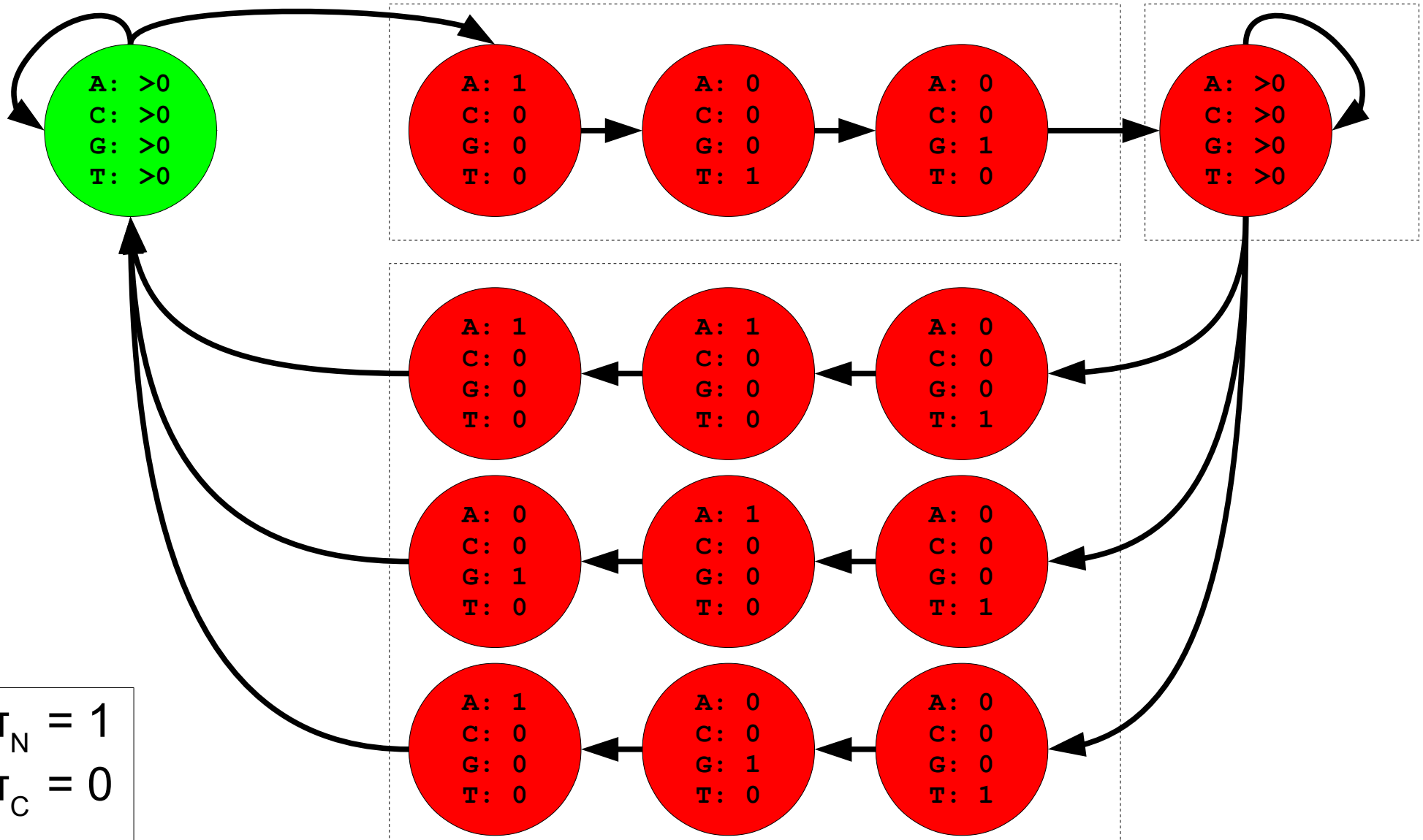


Gene structure

- 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

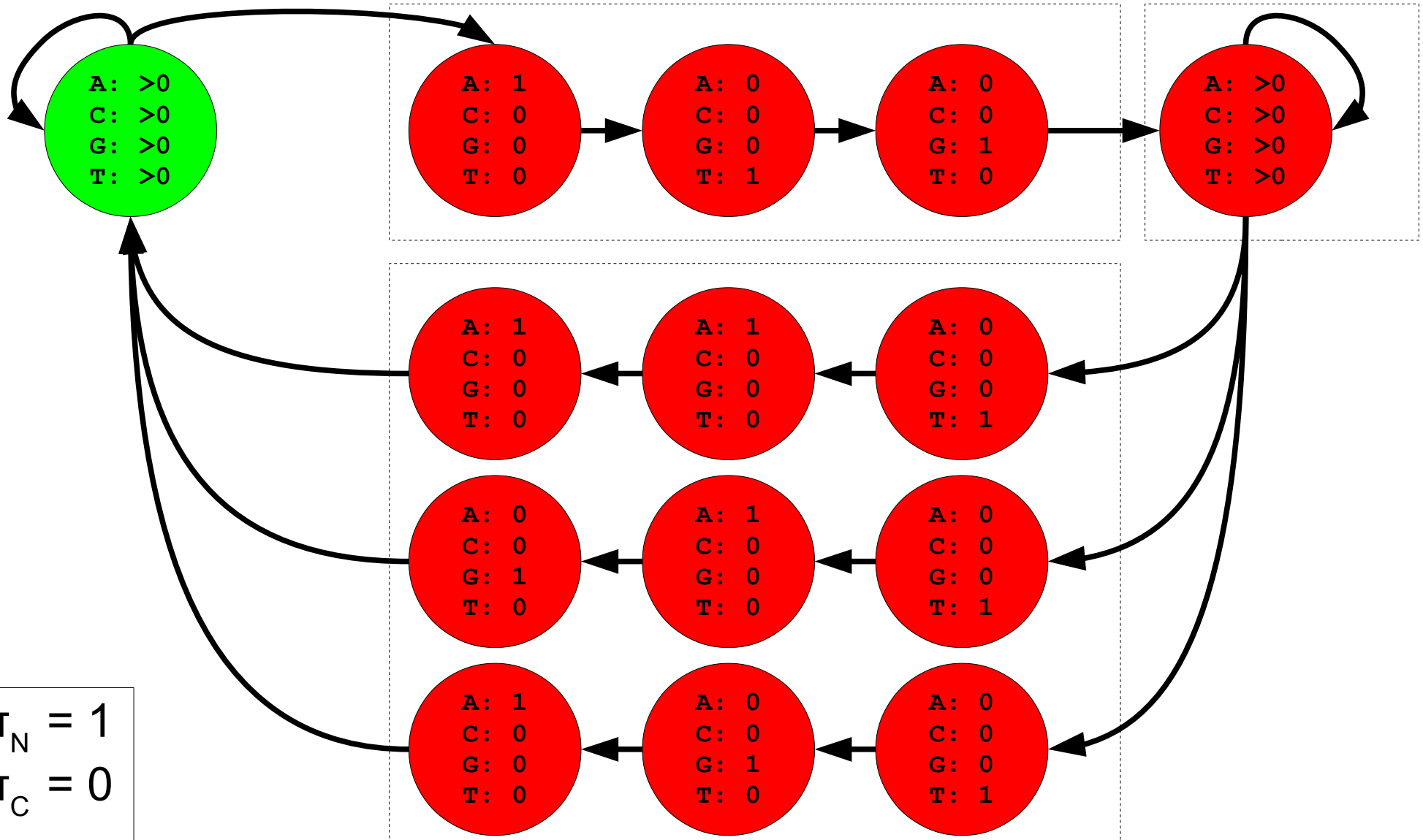


Gene structure

- 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 multiplum of 3

N: non-coding

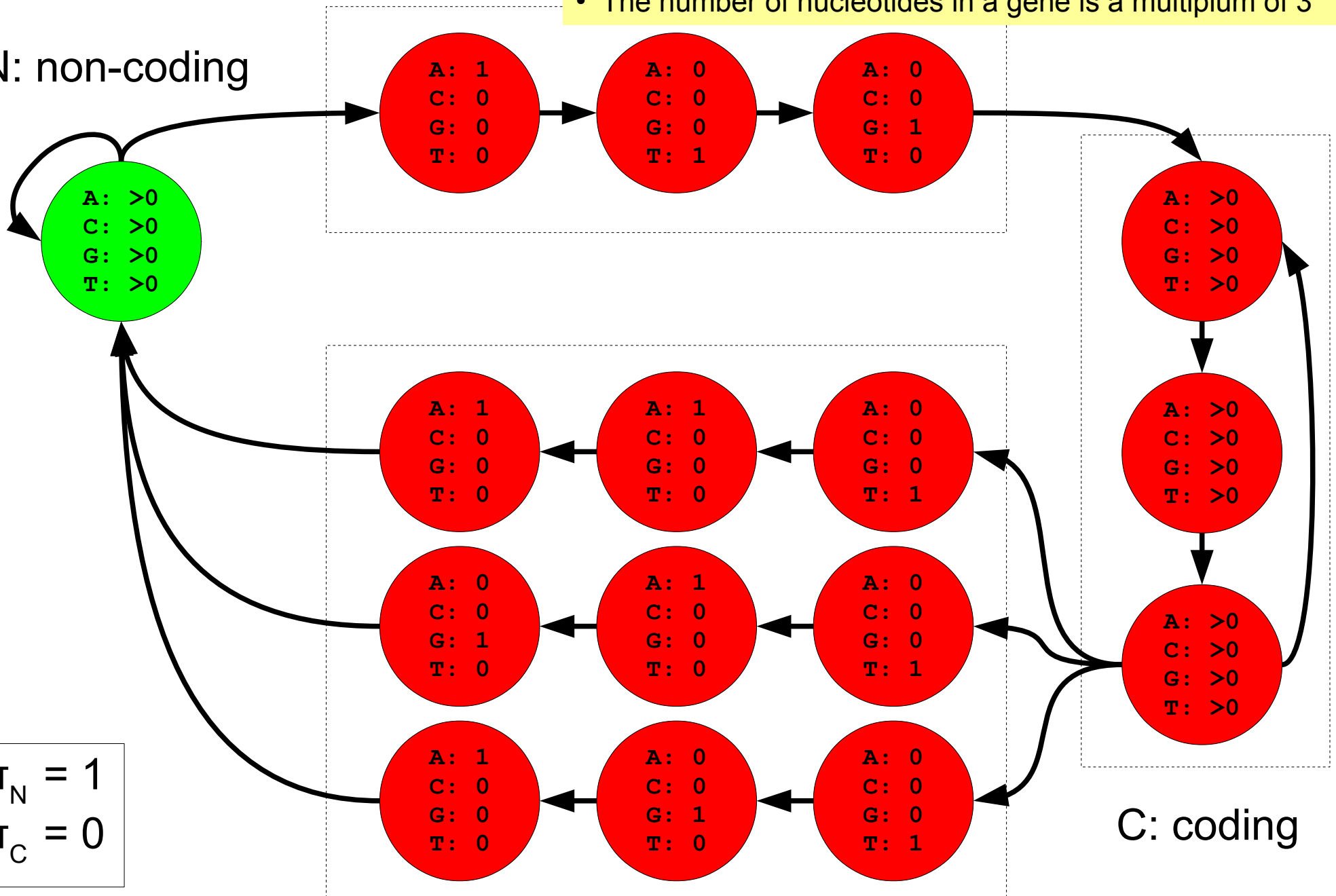
C: coding



Gene structure

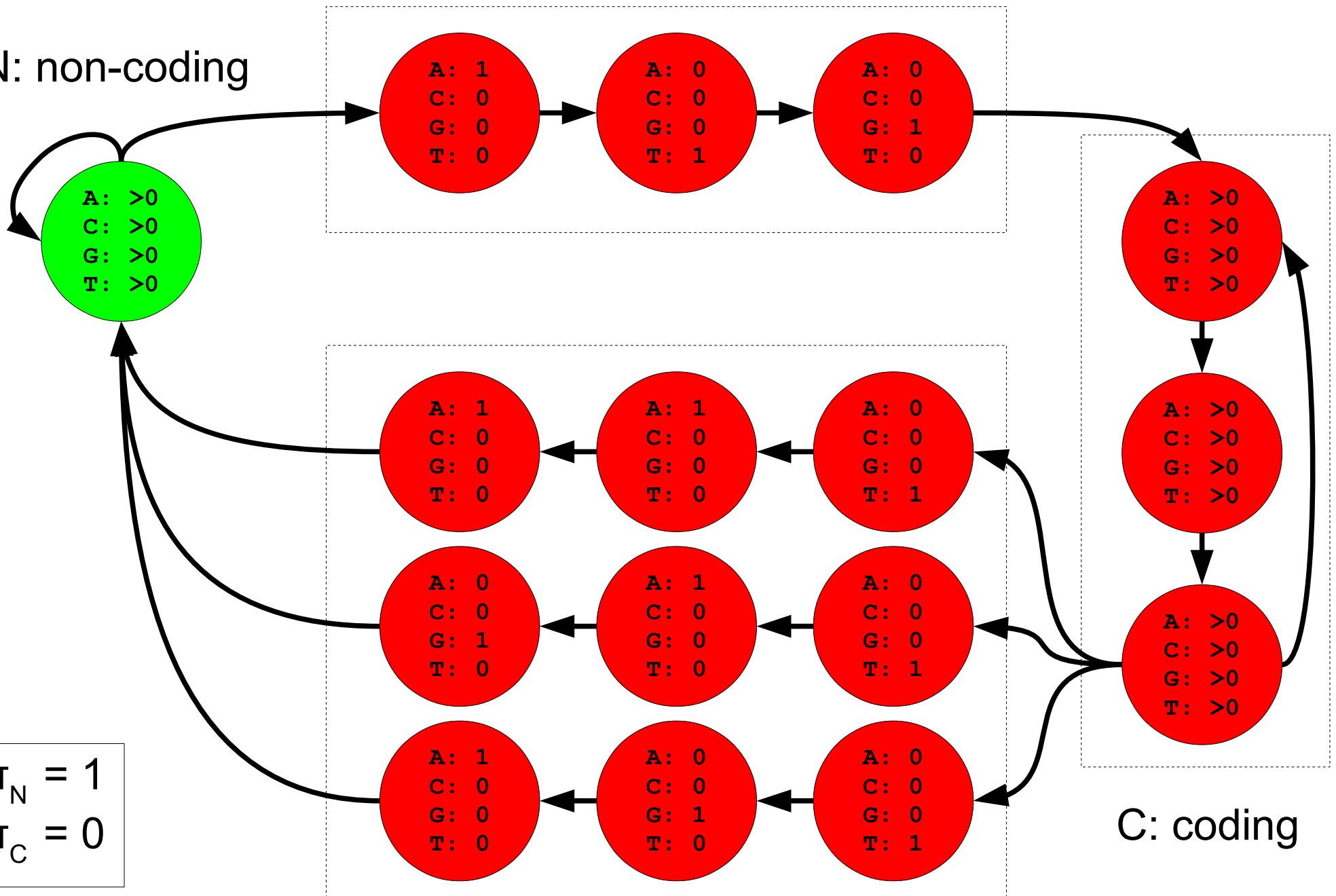
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**
- The gene ends with a stop-codon **taa**, **tag** or **tga**
- The number of nucleotides in a gene is a multiplum of 3

N: non-coding

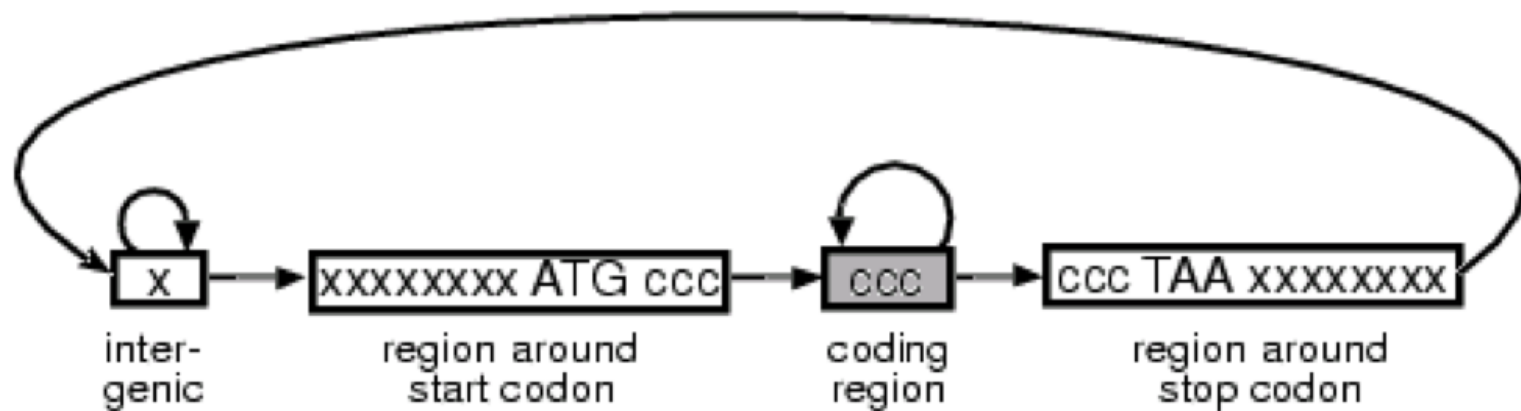


Gene structure in procaryotes

N: non-coding



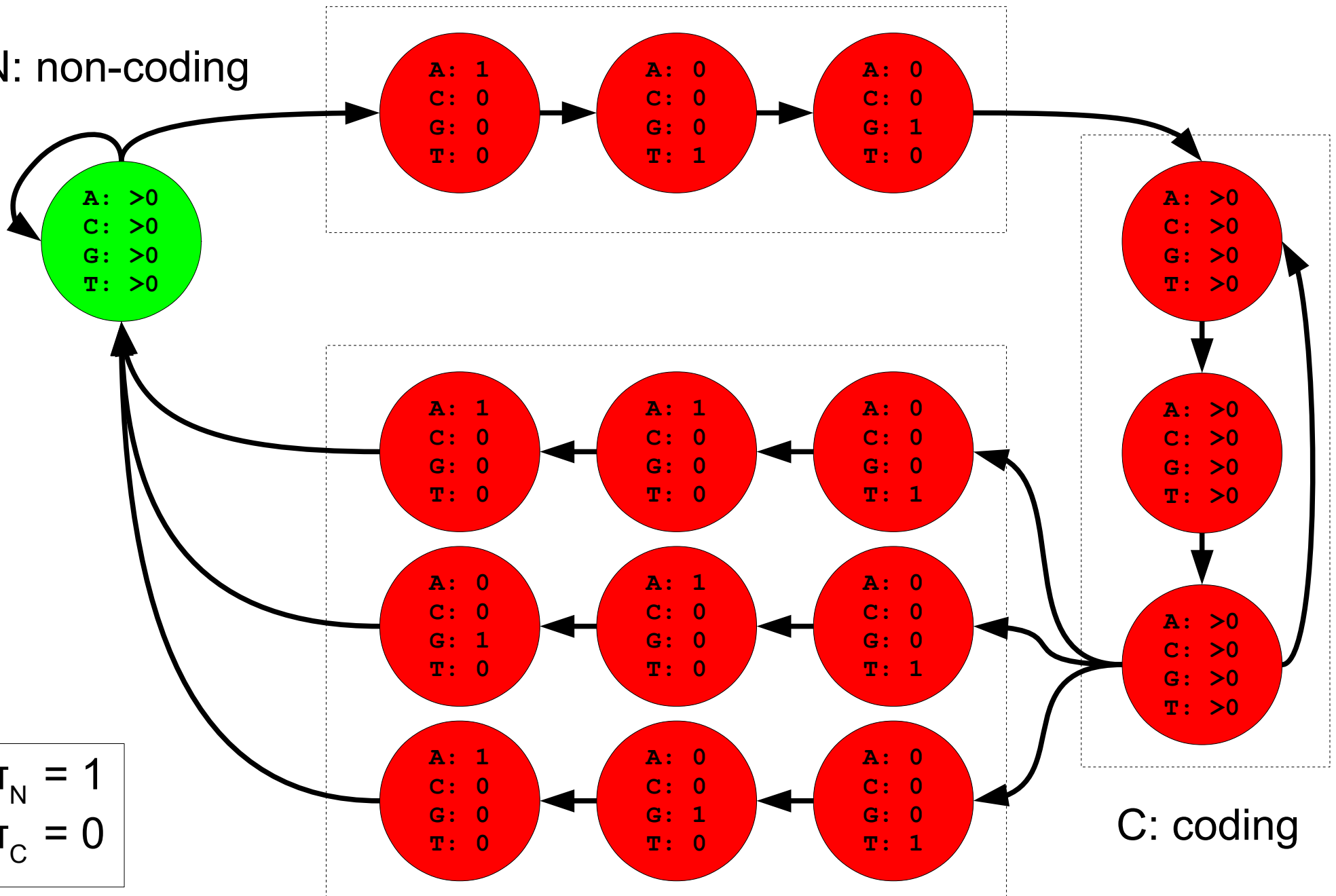
Gene structure in procaryotes



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

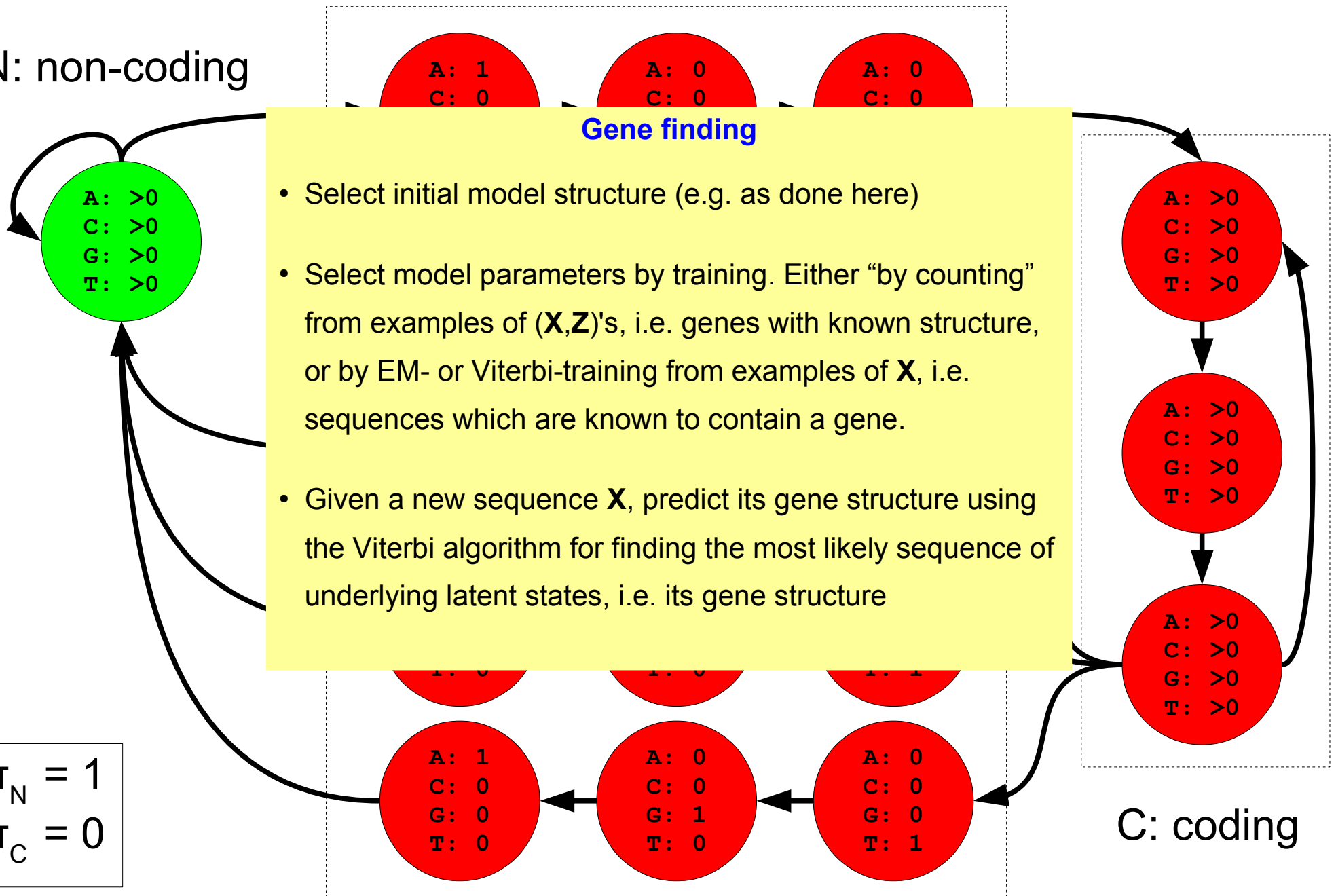
Gene structure in procaryotes

N: non-coding



Gene structure in procaryotes

N: non-coding



Example – Gene finding

N: non-coding

A: 1
C: 0

A: 0
C: 0

A: 0
C: 0

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

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

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

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

C: coding

$$\pi_N = 1$$

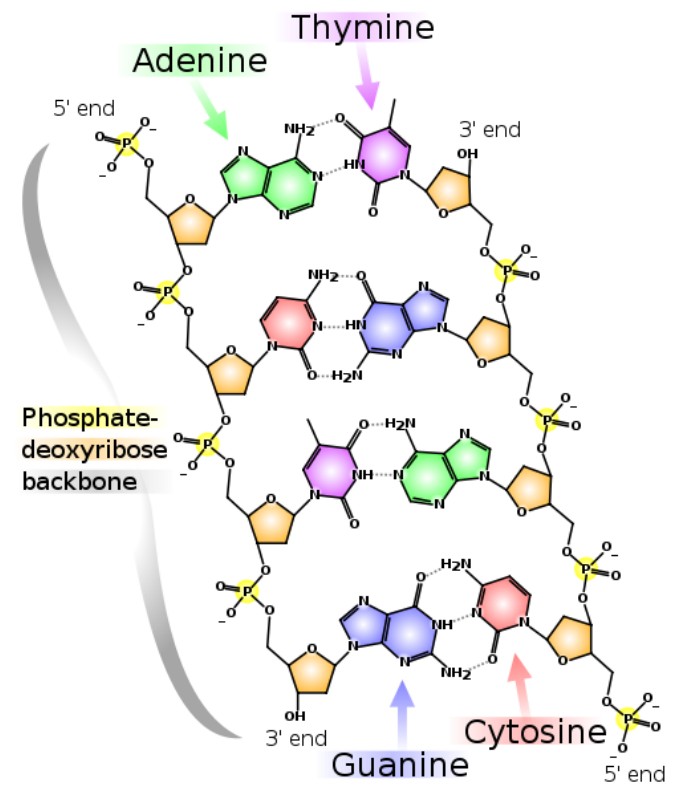
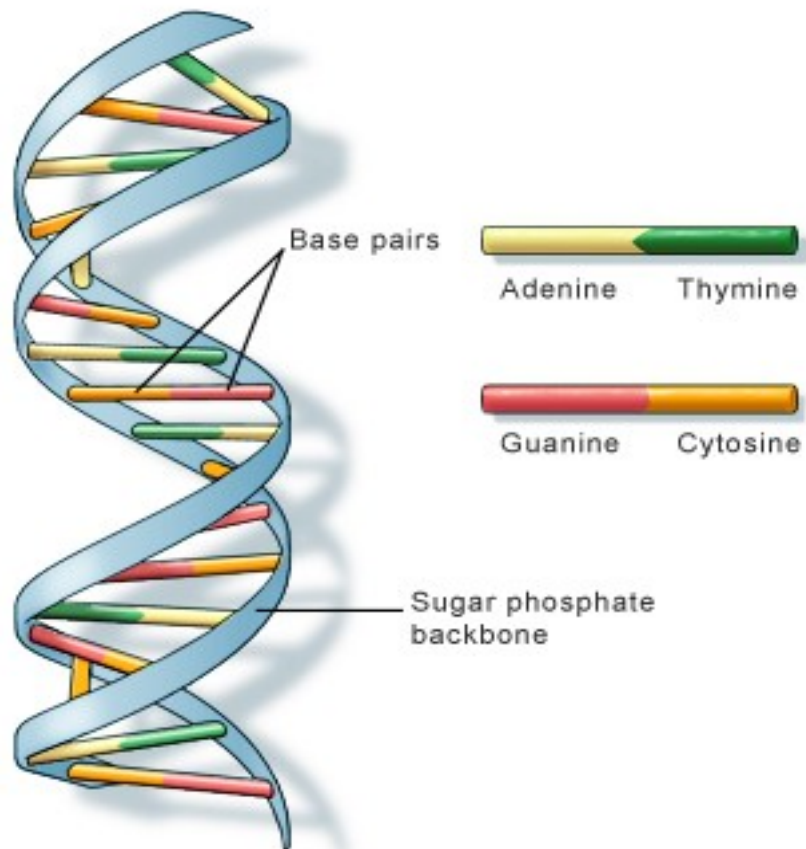
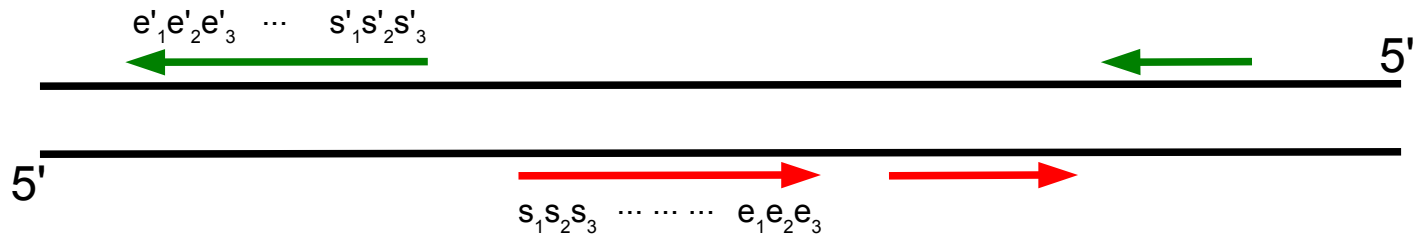
$$\pi_C = 0$$

T: 0

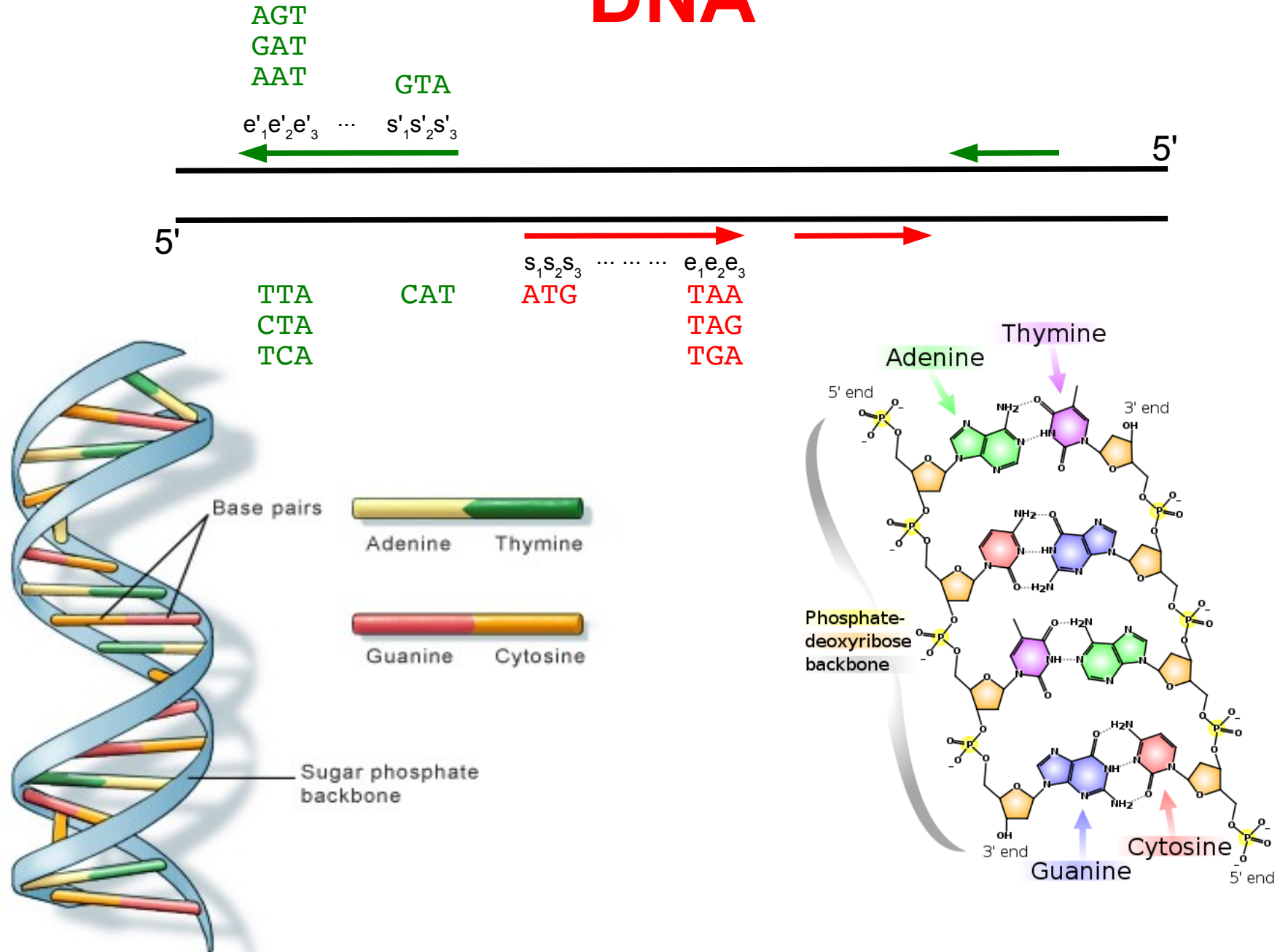
T: 0

T: 1

DNA



DNA



Even more biology

There can be genes in both directions



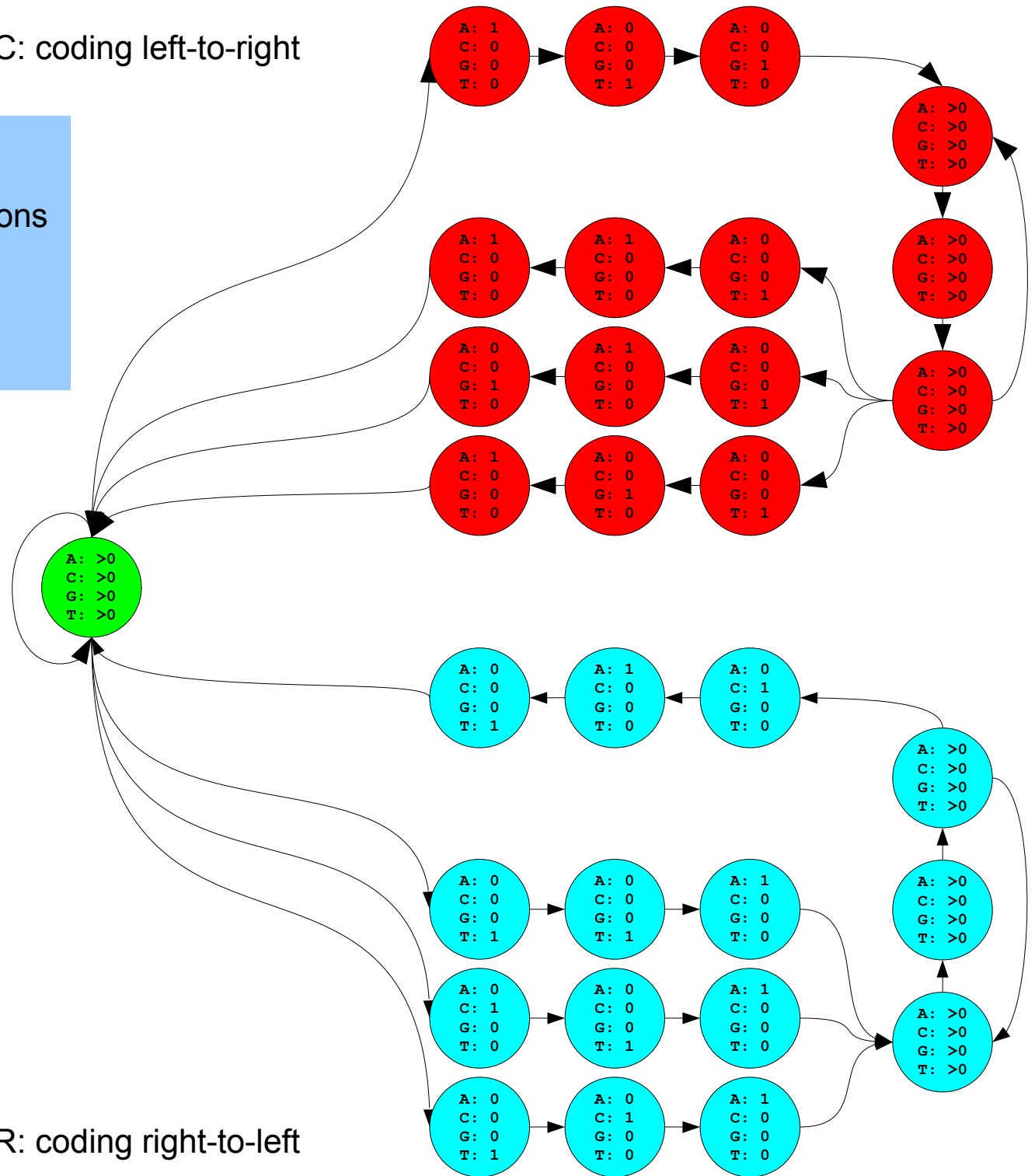
C: coding left-to-right

N: Non-coding

R: coding right-to-left

$$\pi_N = 1$$

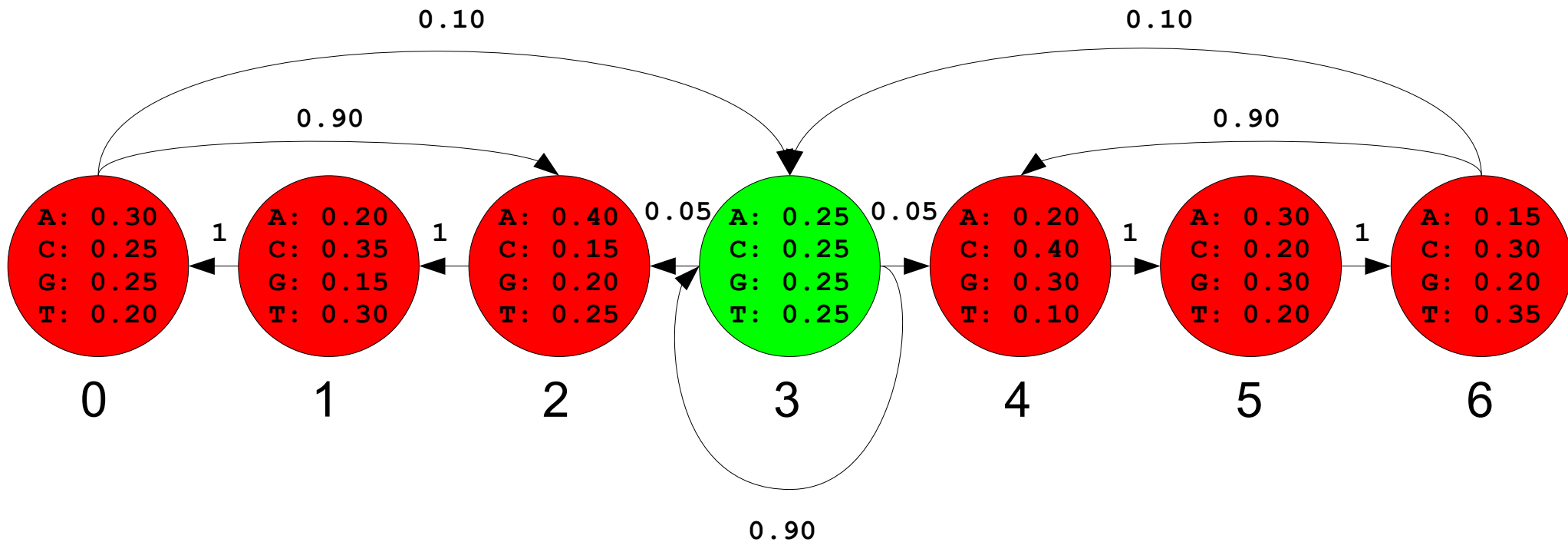
$$\pi_C = 0$$



Example – 7-state HMM

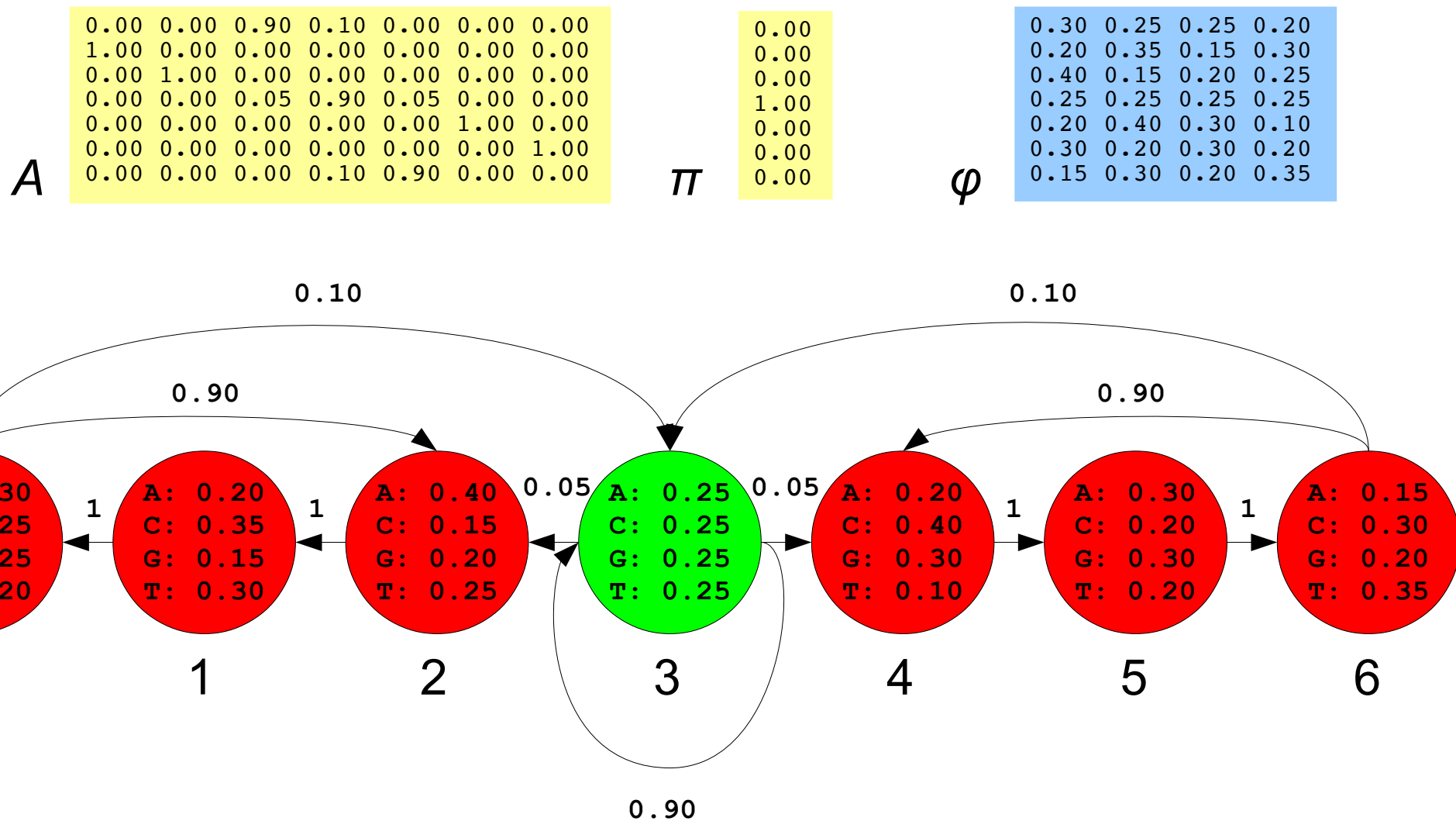
Observable: {A, C, G, T}, States: {0, 1, 2, 3, 4, 5, 6}

A	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	0.00	1.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.00	0.00	0.00	0.00	0.00	0.00	0.00
φ	0.30	0.25	0.25	0.20	0.20	0.30	0.20
	0.20	0.35	0.15	0.30	0.40	0.15	0.30
	0.40	0.15	0.20	0.25	0.25	0.25	0.25
	0.25	0.25	0.25	0.25	0.20	0.40	0.30
	0.20	0.40	0.30	0.10	0.30	0.20	0.30
	0.30	0.20	0.30	0.20	0.15	0.30	0.20
	0.15	0.30	0.20	0.35	0.20	0.20	0.35
	0.15	0.30	0.20	0.35	0.20	0.20	0.35

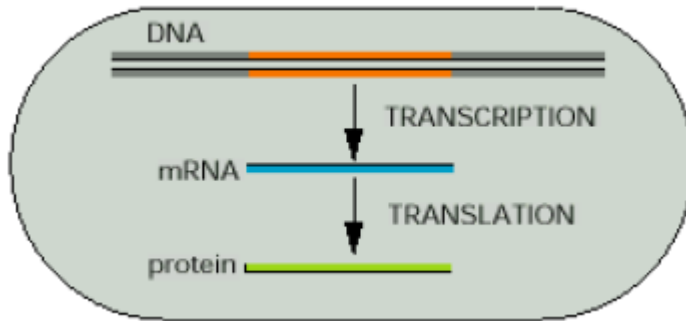


This model is also applicable for gene finding.

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



Problem: From annotation to Z



Biological facts

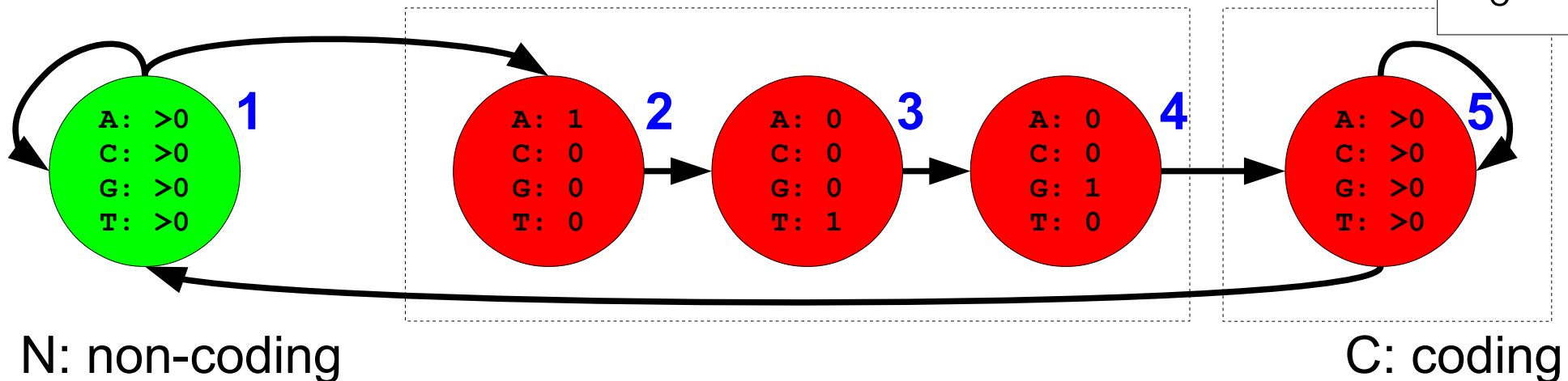
- The gene is a substring of the DNA sequence of A,C,G,T's
- The gene starts with a start-codon **atg**

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaatatgtccgatgacgtgagcataagcgacat

$$\pi_N = 1$$

$$\pi_C = 0$$



Problem: From annotation to Z

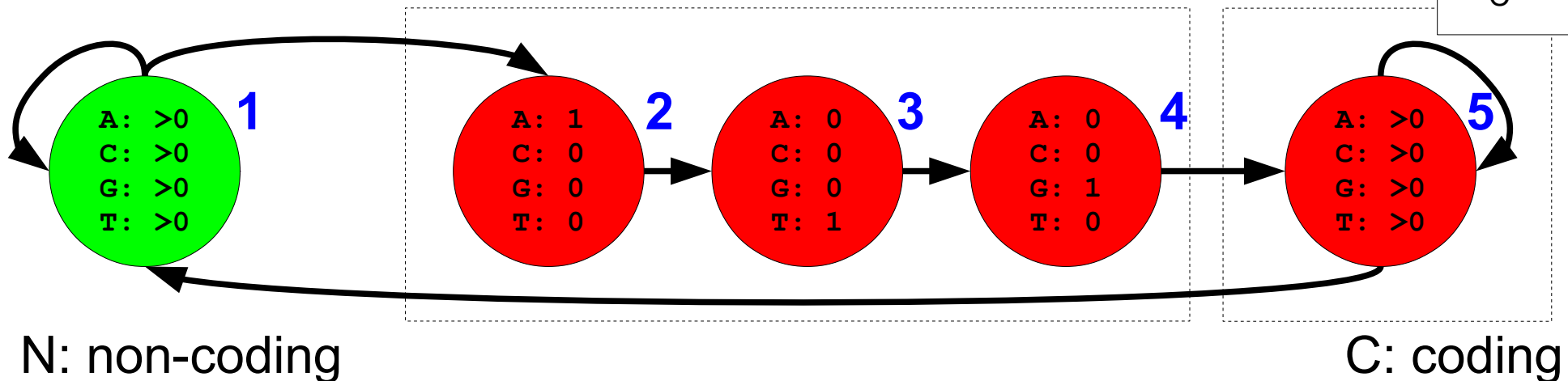
Problem: The string **Z**=NNNCCC.... is not a proper sequence of states in the illustrated HMM, but it can easily be converted into one (because there is this case is a 1-1 matching between a sequence of Ns and Cs and a sequence of states).

ence of A,C,G,T's

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

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



Problem: From annotation to Z

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

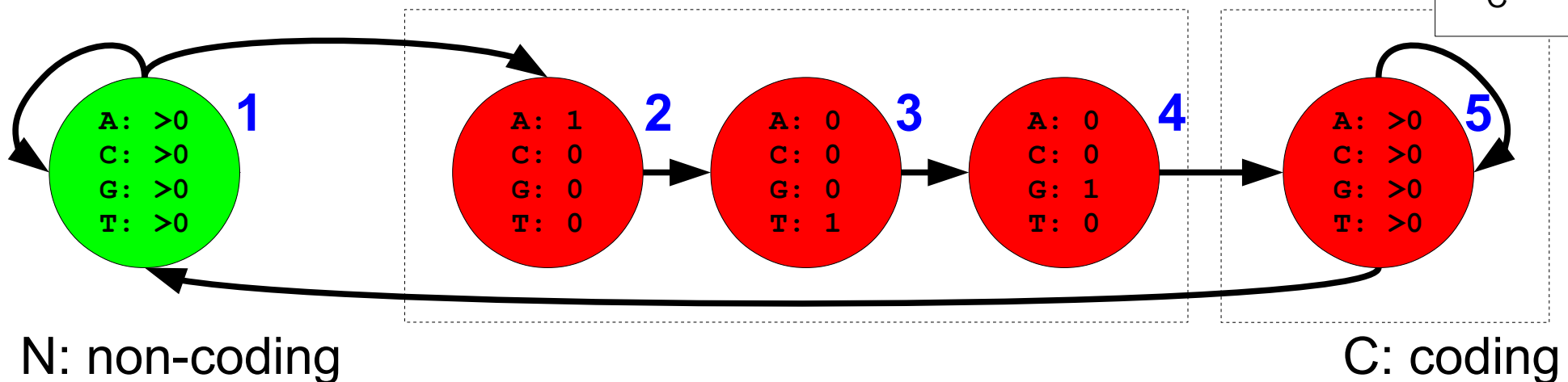
ence of A,C,G,T's

1112345555551111111123455555555555555111111111111

Z: NNNCCCCCCCCNNNNNNNNCCCCCCCCCCCCCCCCNNNNNNNNNNNN

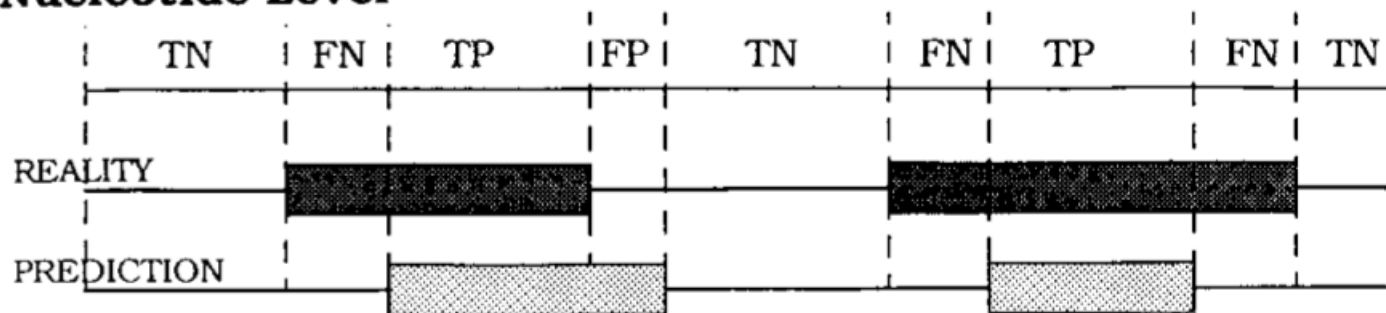
X: acgatgcgctaataatgtccgatgacgtgagcataagcgacat

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



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

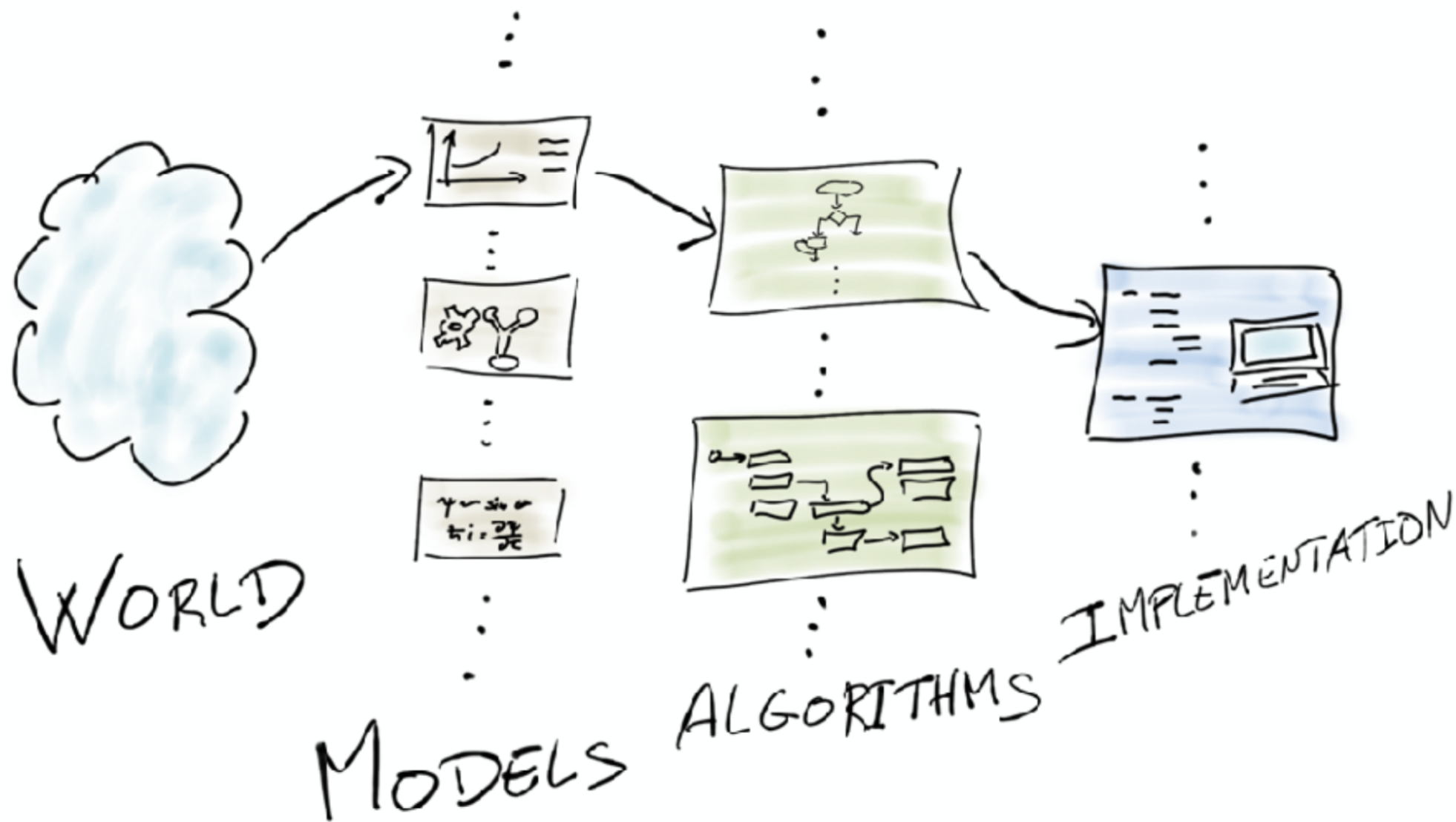
$$CC = \frac{(TP \times TN) - (FN \times FP)}{\sqrt{(TP + FN) \times (TN + FP) \times (TP + FP) \times (TN + FN)}}$$

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



Which model should I choose?

Relevance: Which model do I believe capture “real life” best?

Applicable: Which model can be used in practice?

Do my computer have memory enough to handle the computation?

Is my computer “fast enough” to perform the computation?

What does “fast enough” mean, and what influences running time?

How would you answer these questions?

Exercise

Consider the following simple “program” that has the same “algorithmic complexity” as the Viterbi algorithm:

```
sum = 0
for n = 1 to N:
    for k = 1 to K:
        for j = 1 to K:
            sum = sum + 1
print sum
```

Try to implement it in different programming languages that you know and see how long time it takes for realistic choices of N and K, e.g. $N \approx 2.000.000$ and $K = 7$.