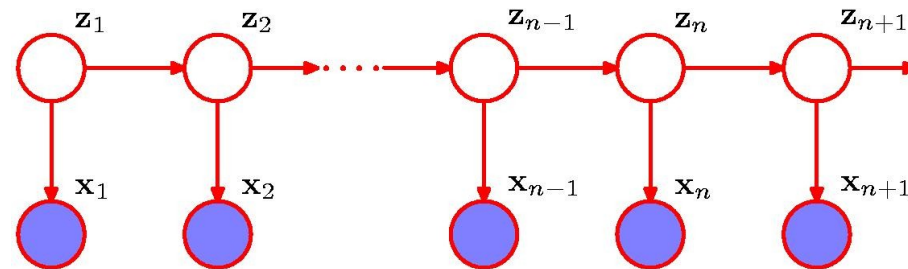


Hidden Markov Models

Some useful extensions



Christian Nørgaard Storm Pedersen

cstorm@birc.au.dk

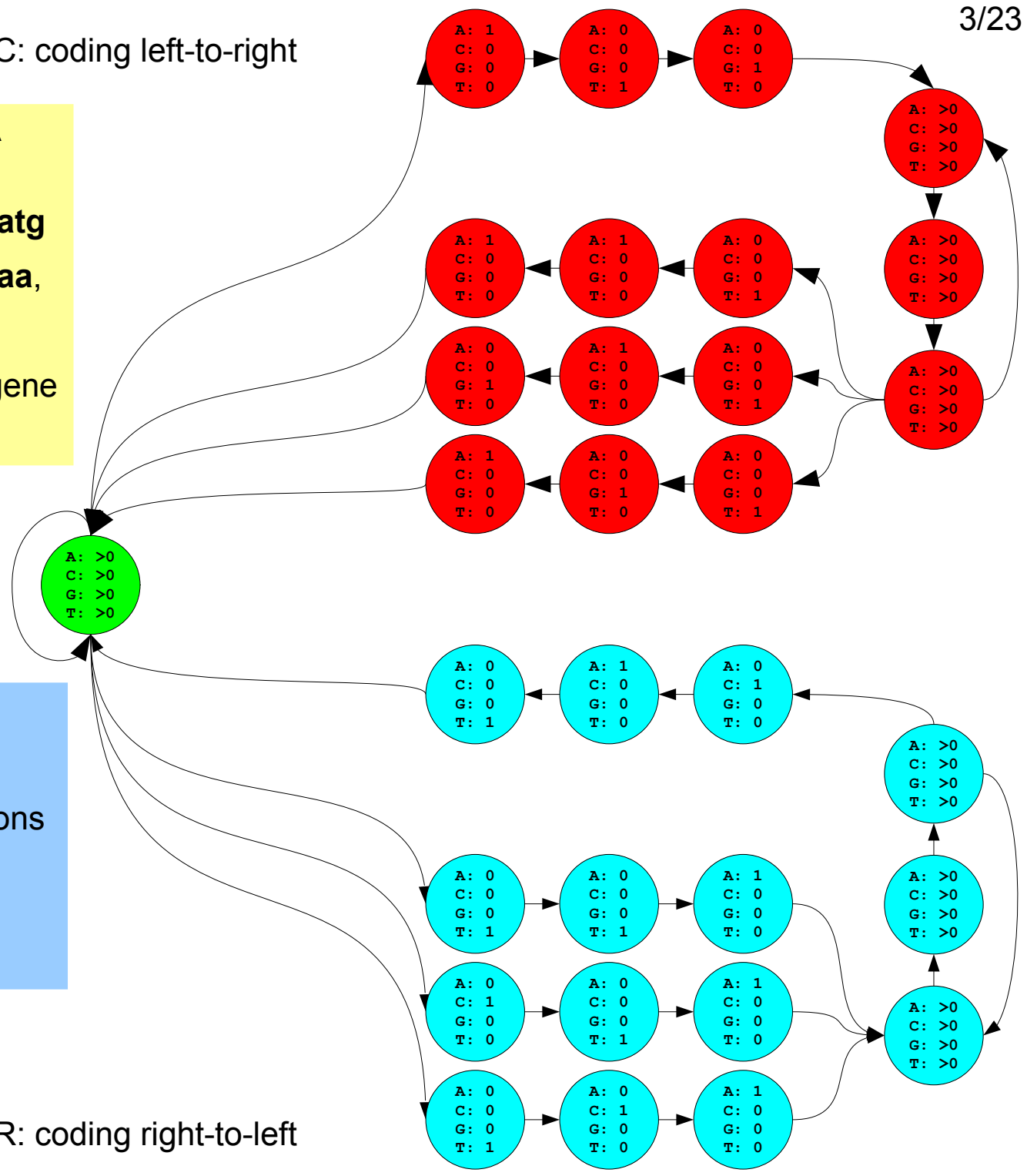
Last time

We developed an HMM for gene finding ...

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

C: coding left-to-right

N: Non-coding



Even more biology

There can be genes in both directions

$$\pi_N = 1$$

$$\pi_C = 0$$

R: coding right-to-left

C: coding left-to-right

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

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

N: No

Even more b

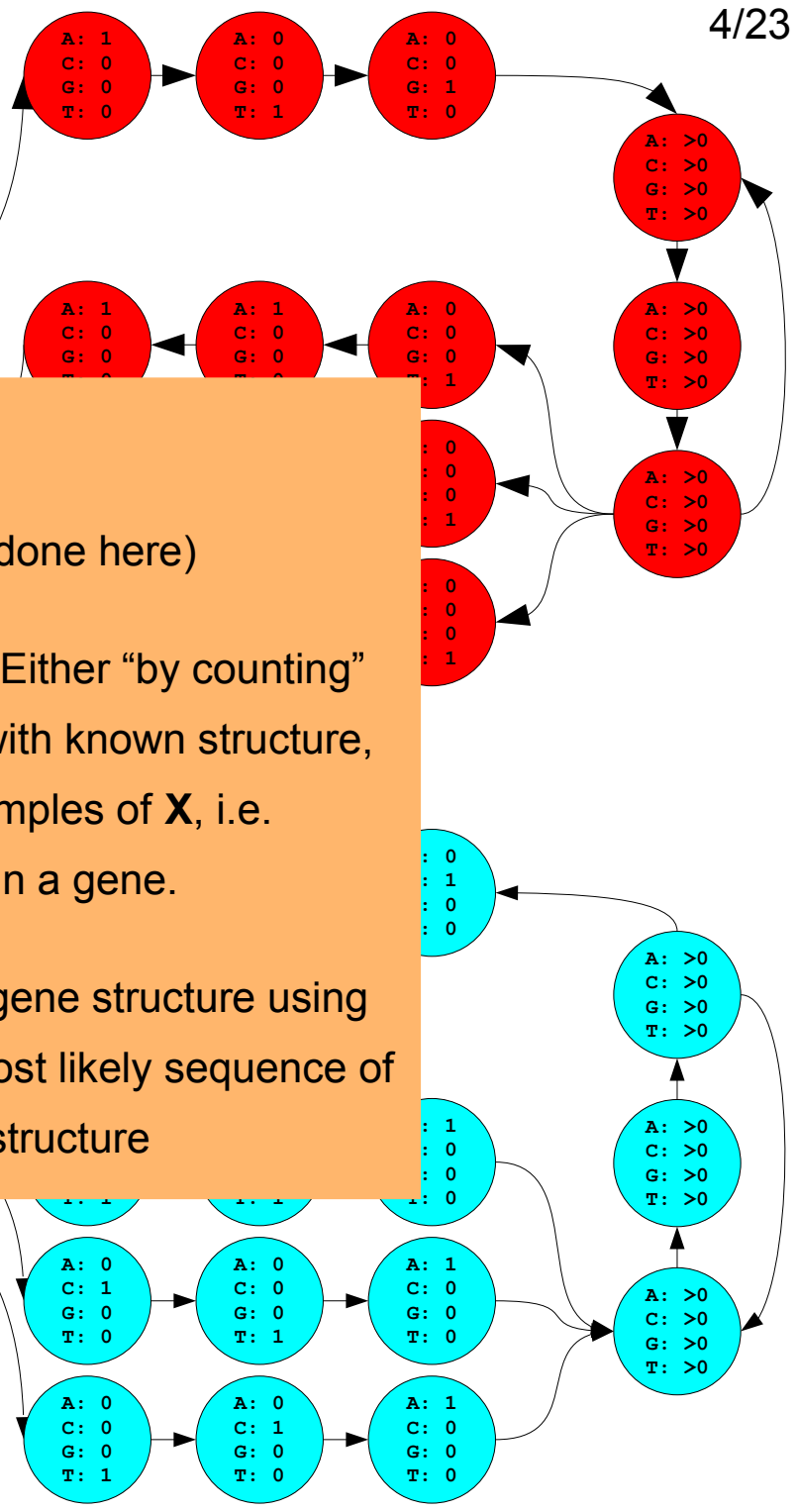
There can be genes in



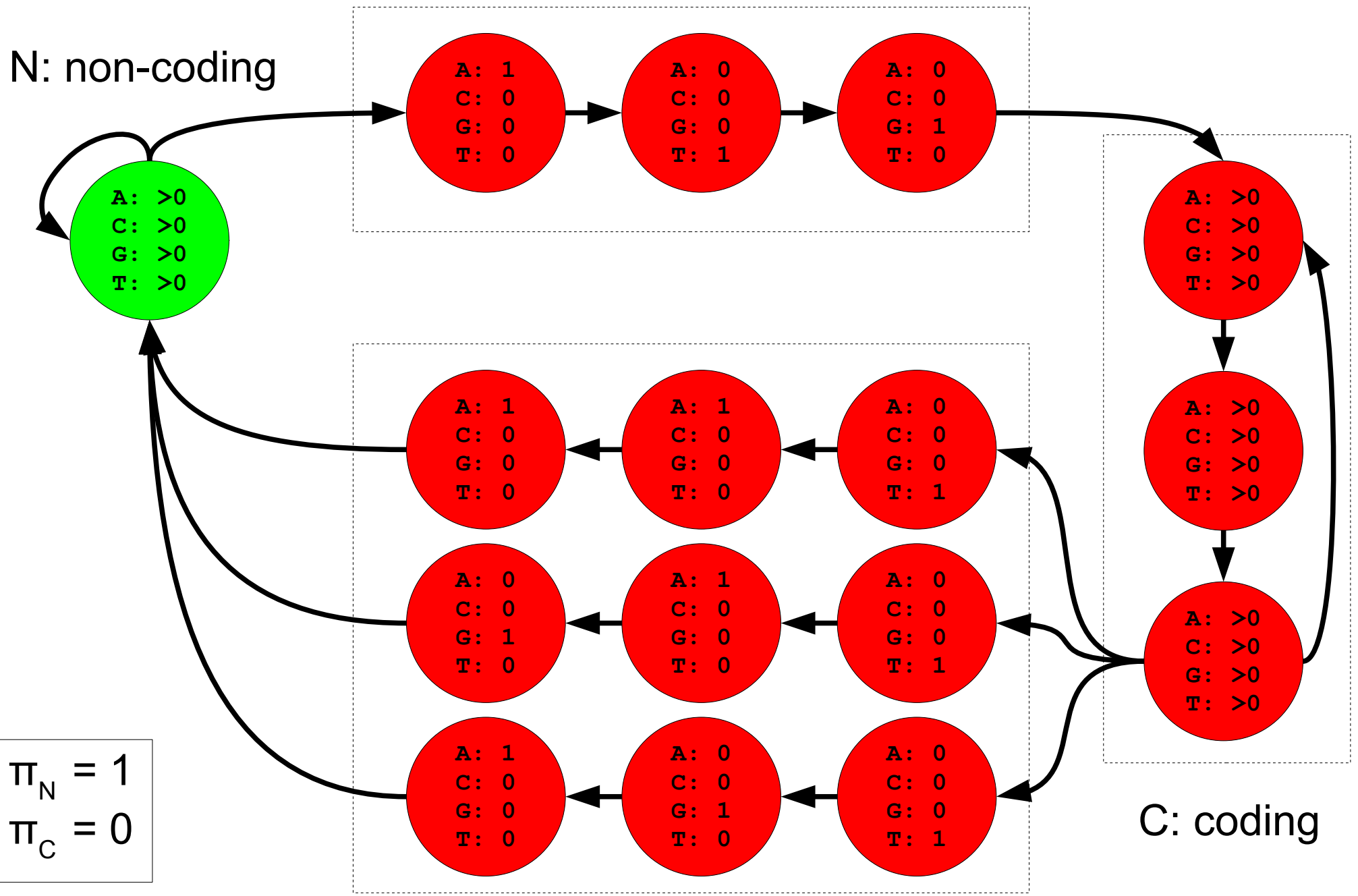
$$\pi_N = 1$$

$$\pi_C = 0$$

R: coding right-to-left



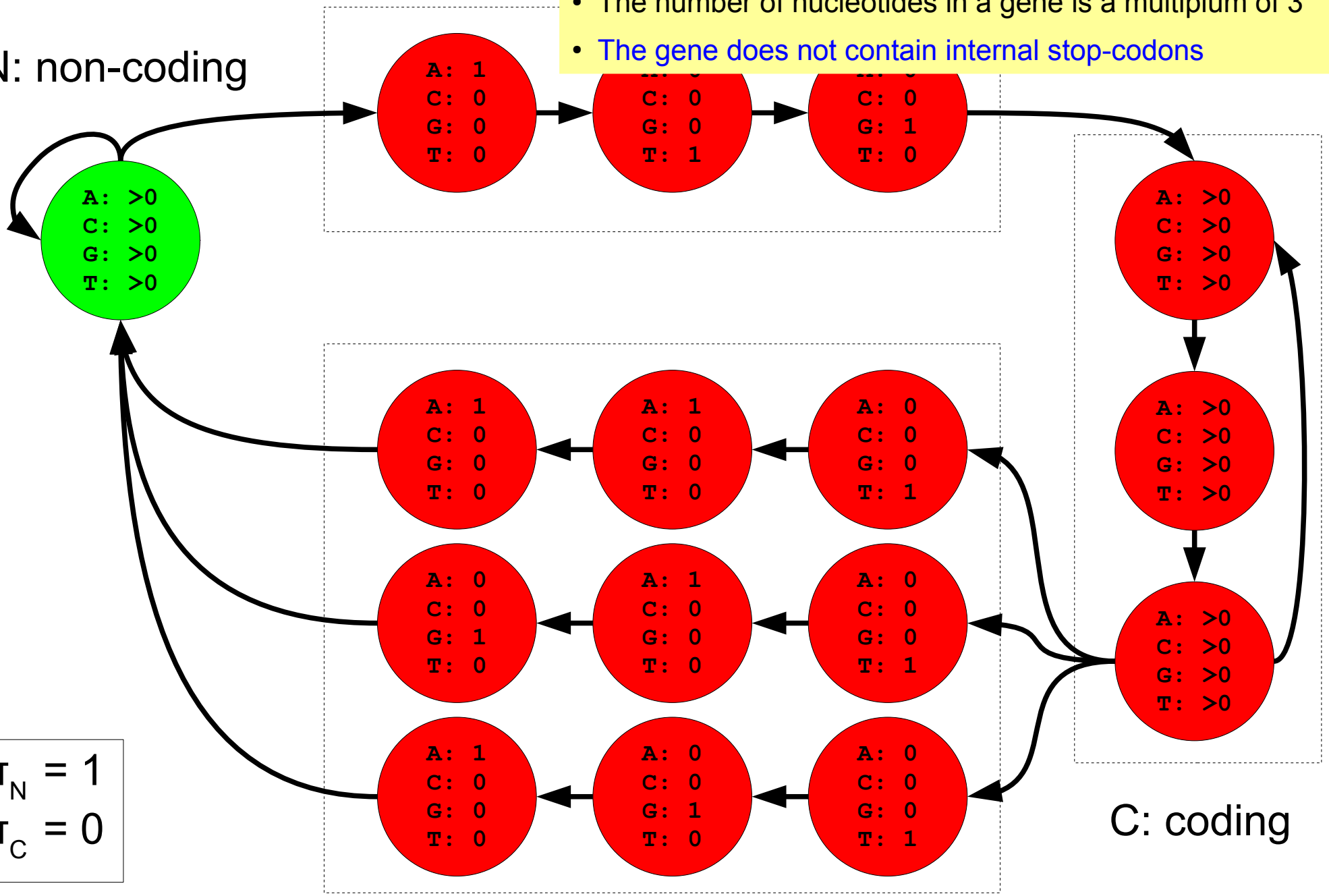
The "forward-coding" part



The "forward"

- 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
- The gene does not contain internal stop-codons

N: non-coding



$$\pi_N = 1$$

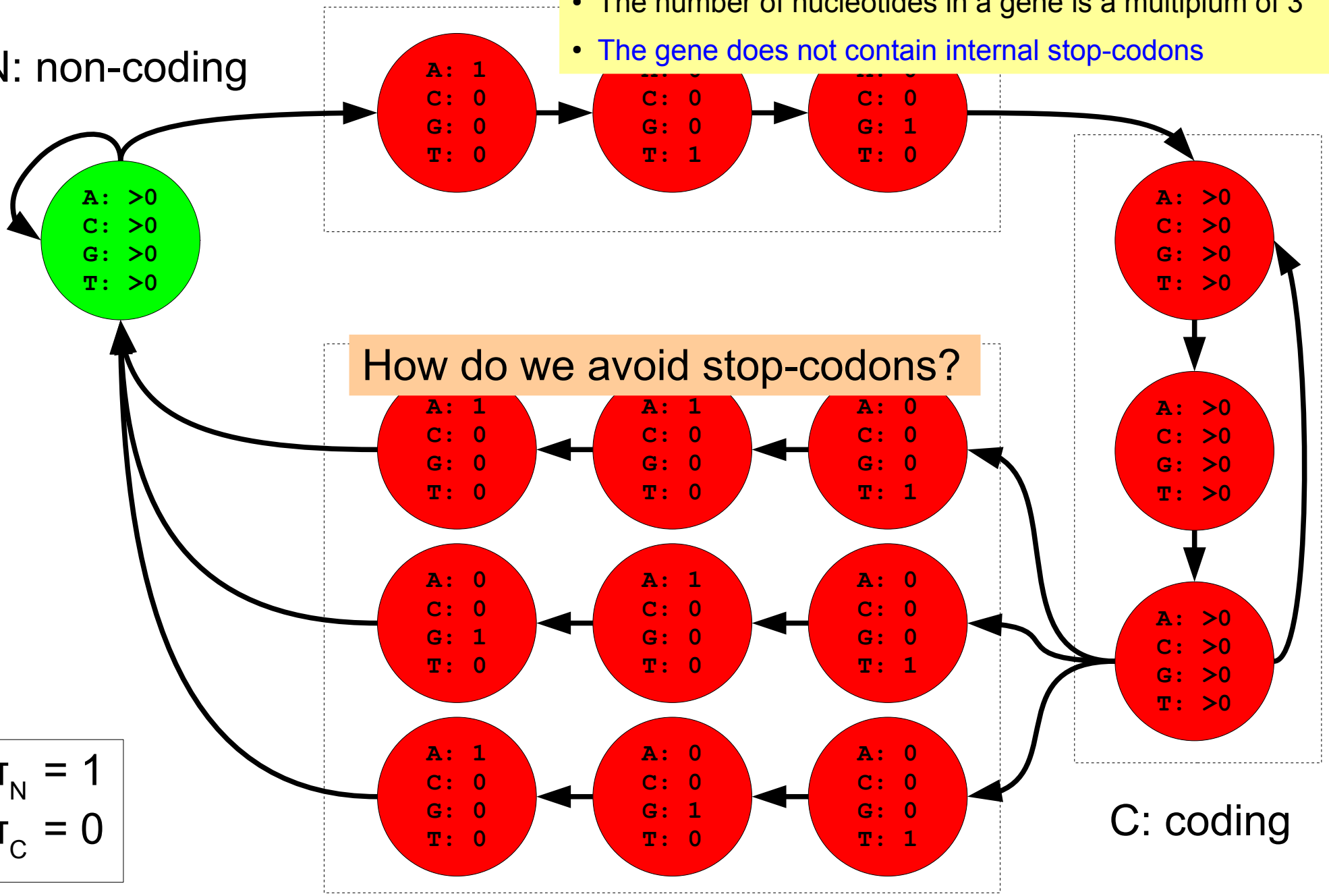
$$\pi_C = 0$$

C: coding

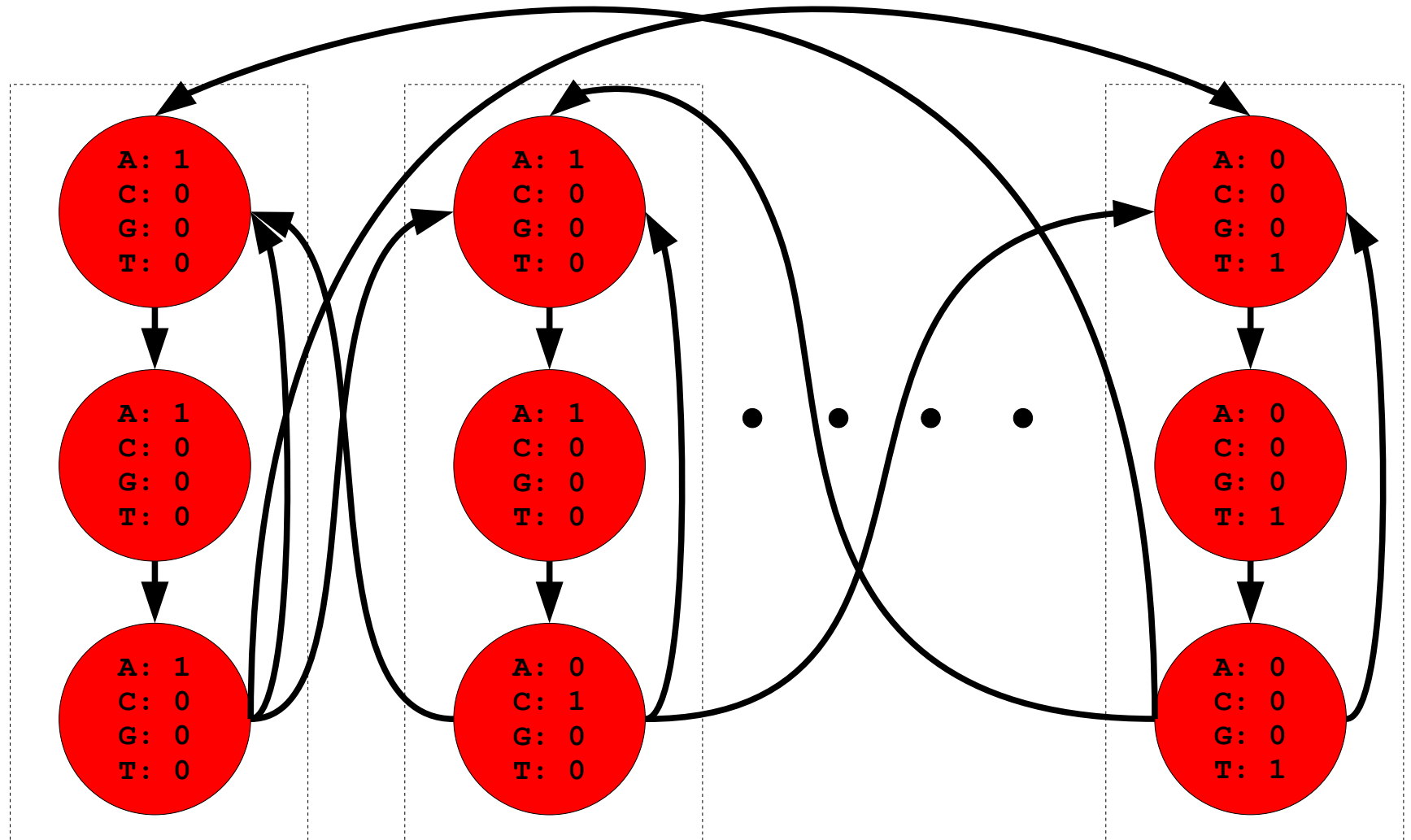
The "forward"

- 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
- The gene does not contain internal stop-codons

N: non-coding



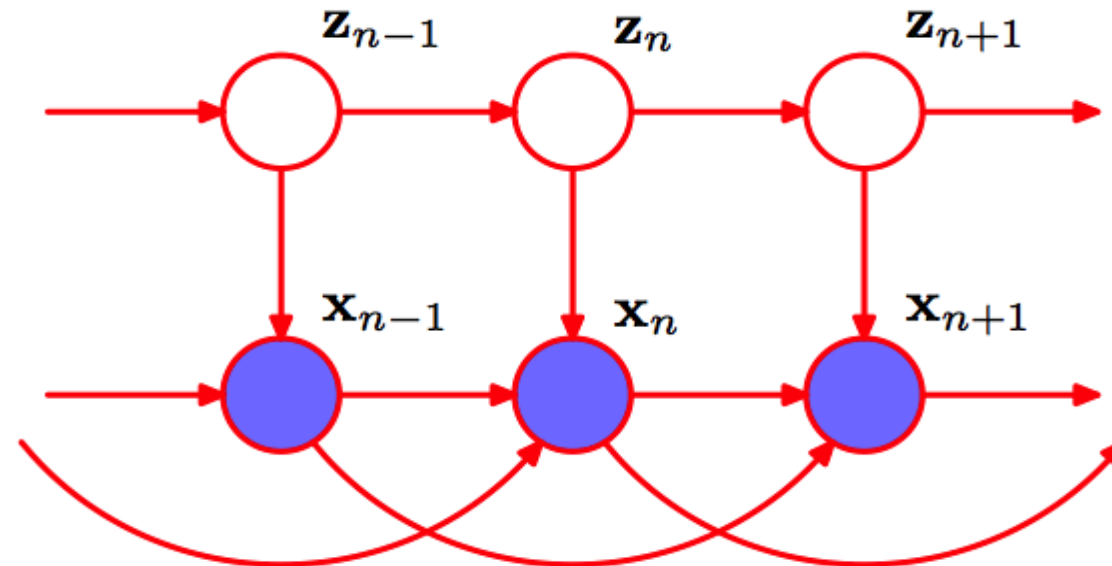
Avoiding internal stop-codons



Encode the emission of each legal codon as a sequence of states.
 Many states ($61 \cdot 3 = 183$) and non-trivial transitions ($61 \cdot 60 = 3660$)!

Other ideas?

Autoregressive HMMs

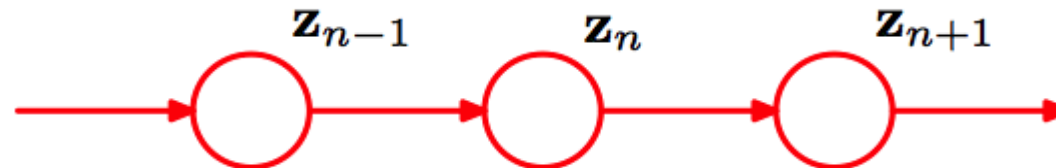


The probability of emitting \mathbf{x}_n depends also on \mathbf{x}_{n-1} and \mathbf{x}_{n-2}
 The basic algorithms remain the same:

$$\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

$$\omega(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, \mathbf{z}_n) \max_{\mathbf{z}_{n-1}} \omega(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

Autoregressive HMMs



For each state, we just have to state the conditional probabilities. For a 4-letter DNA alphabet this corresponds to $4 \times 4 \times 4$ emission prob.



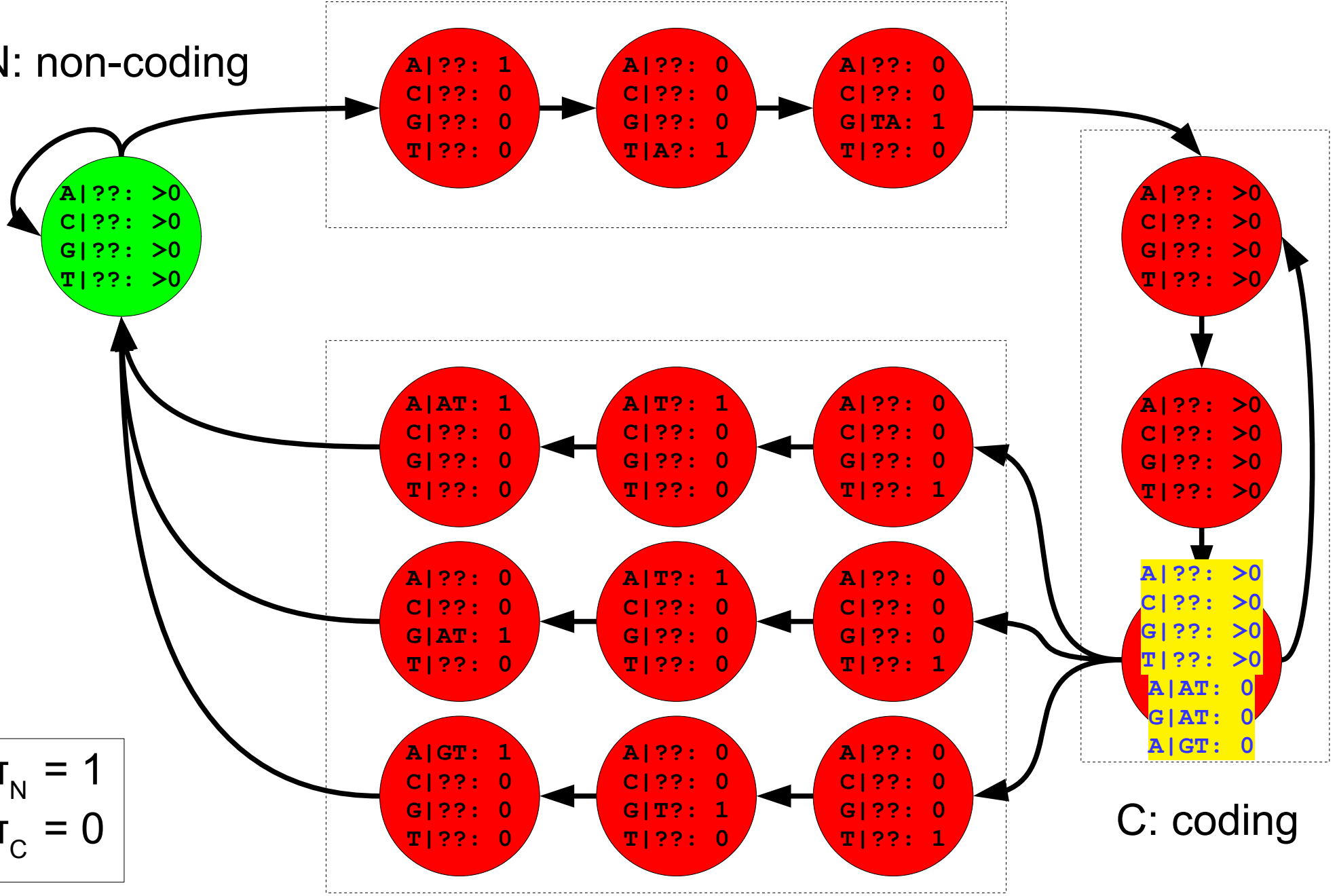
The probability of emitting \mathbf{x}_n depends also on \mathbf{x}_{n-1} and \mathbf{x}_{n-2}
The basic algorithms remain the same:

$$\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

$$\omega(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, \mathbf{z}_n) \max_{\mathbf{z}_{n-1}} \omega(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

Adjusting our simple HMM

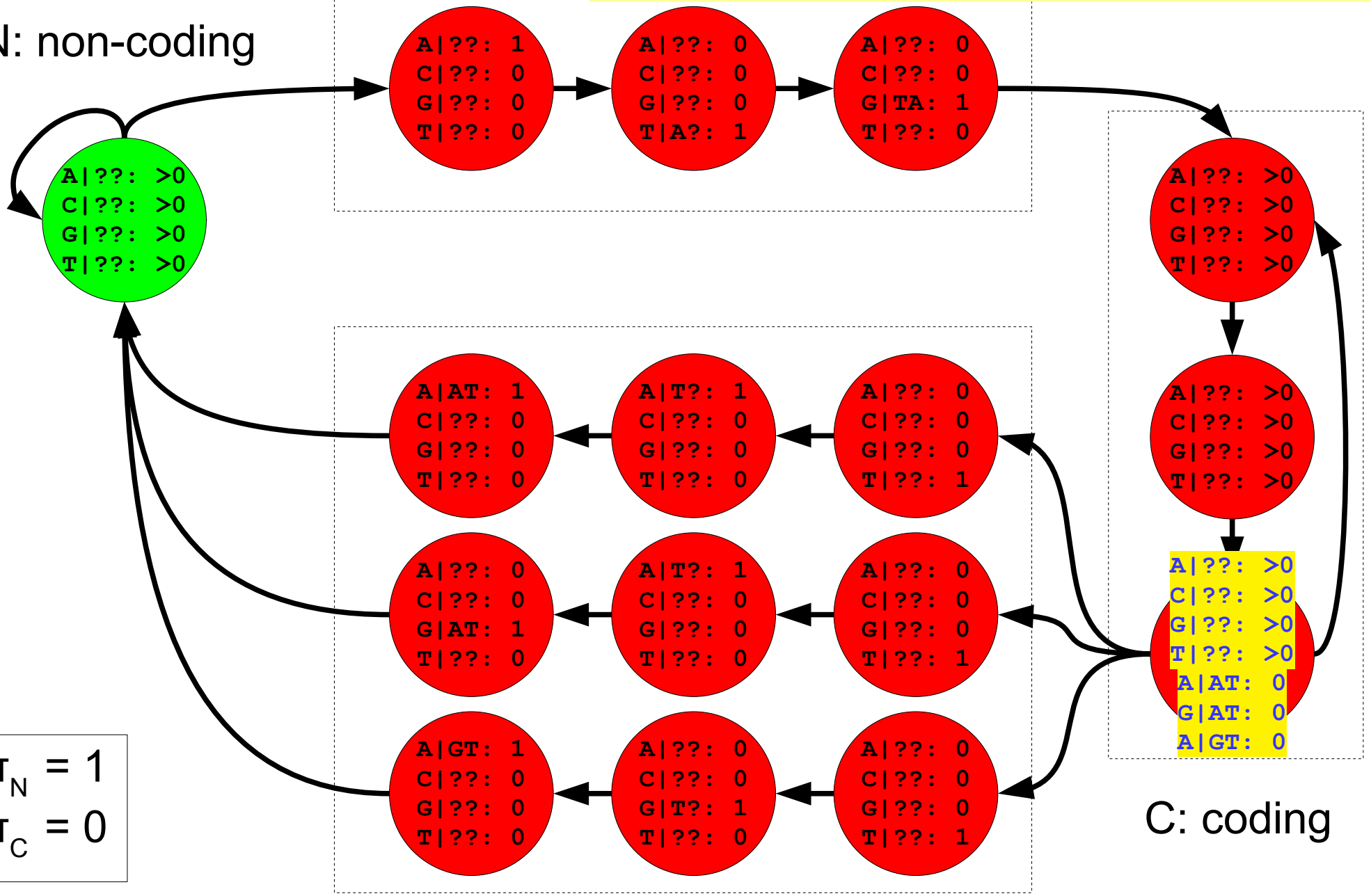
N: non-coding



- 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 gene does not contain internal stop-codons

Adjusting

N: non-coding



$$\pi_N = 1$$

$$\pi_C = 0$$

C: coding

Emitting a variable number of symbols

Make it possible to emit a variable number of symbols depending on the state. Fx when being in state \mathbf{z}_n the model emits $d_{\mathbf{z}_n}$ symbols, where $d_{\mathbf{z}_n}$ is an integer ≥ 0 .

The basic algorithms can easily be reformulated, fx Viterbi:

Emitting a variable number of symbols

Make it possible to emit a variable number of symbols depending on the state. Fx when being in state \mathbf{z}_n the model emits $d_{\mathbf{z}_n}$ symbols, where $d_{\mathbf{z}_n}$ is an integer ≥ 0 .

The basic algorithms can easily be reformulated, fx Viterbi:

$\omega(n, k)$: The probability of the most likely path generating the first n symbols and ending in state k .

$$\omega(n, k) = \max_{k' \rightarrow k} \omega(n - d_k, k') p(k' \rightarrow k) p(\mathbf{x}_n \dots \mathbf{x}_{n-d_k+1} | k)$$

Emitting a variable number of symbols

Make it possible to emit a variable number of symbols depending on the state. Fx when being in state \mathbf{z}_n the model emits $d_{\mathbf{z}_n}$ symbols, where $d_{\mathbf{z}_n}$ is an integer ≥ 0 .

The basic algorithms can easily be reformulated, fx Viterbi:

$\omega(n, k)$: The probability of the most likely path generating the first n symbols and ending in state k .

$$\omega(n, k) = \max_{k' \rightarrow k} \omega(n - d_k, k') p(k' \rightarrow k) p(\mathbf{x}_n \dots \mathbf{x}_{n-d_k+1} | k)$$

Transition prob from state k' to k

Emission prob of emitting d_k symbols from state k .

Emitting a variable number of symbols

Make it possible to emit a variable number of symbols depending on the state. Fx when being in state \mathbf{z}_n the model emits $d_{\mathbf{z}_n}$ symbols, where $d_{\mathbf{z}_n}$ is an integer ≥ 0 .

The basic algorithms can easily be reformulated, fx Viterbi:

$\omega(n, k)$: The probability of the most likely path for the first n symbols and ending in state k .

Special case: If $d_k = 0$ then the state is called a *silent state*.

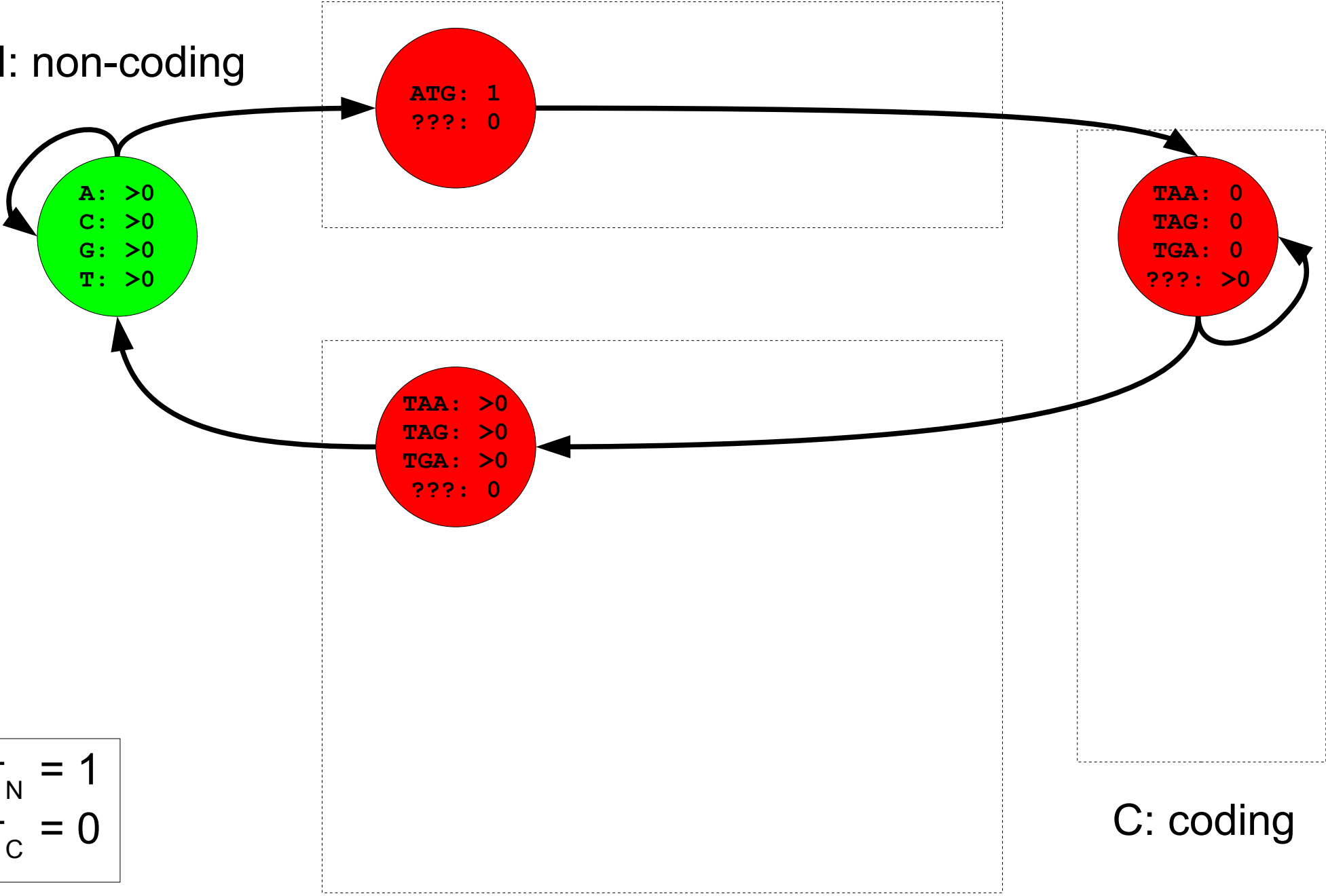
$$\omega(n, k) = \max_{k' \rightarrow k} \omega(n - d_k, k') p(k' \rightarrow k) p(\mathbf{x}_n \dots \mathbf{x}_{n-d_k+1} | k)$$

Transition prob from state k' to k

Emission prob of emitting d_k symbols from state k .

Adjusting our simple HMM

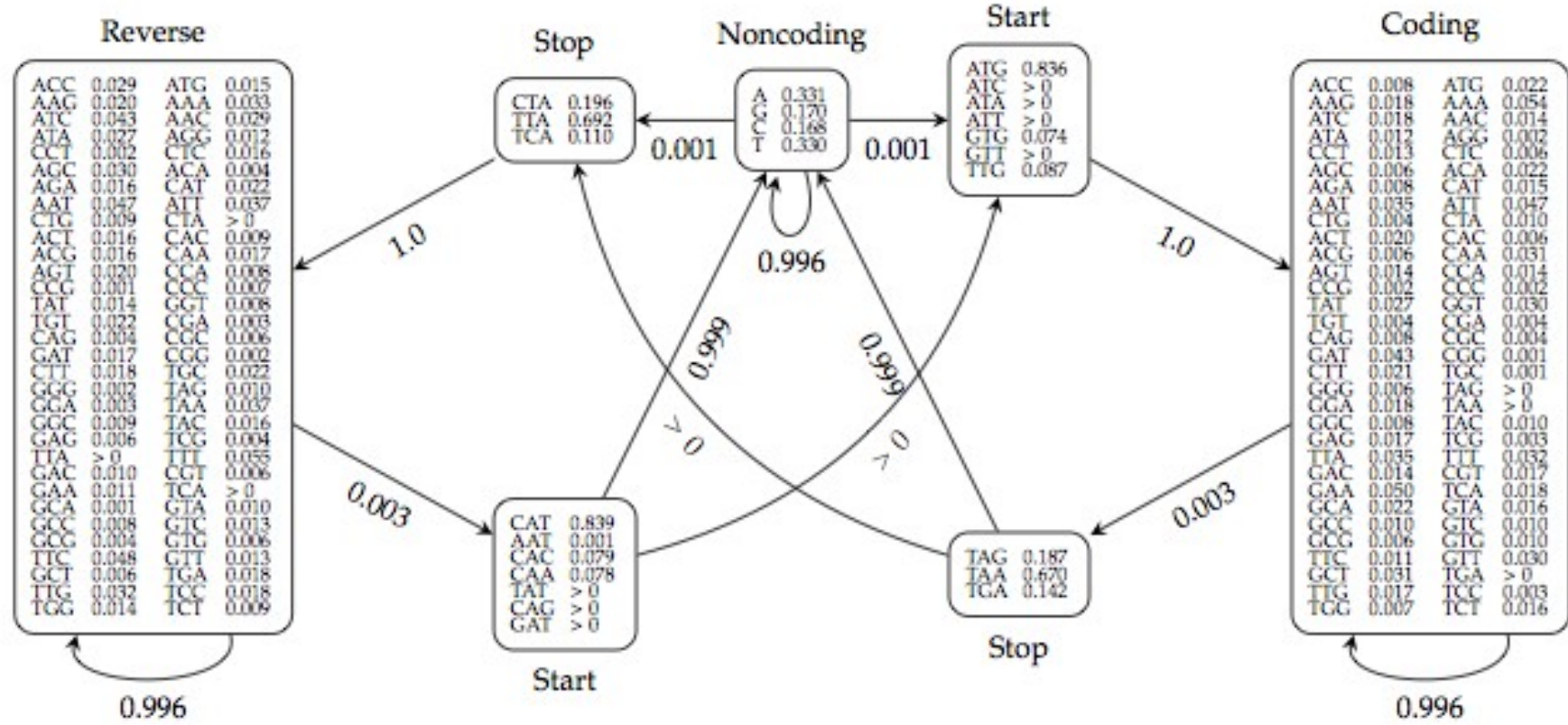
N: non-coding



$$\pi_N = 1$$

$$\pi_C = 0$$

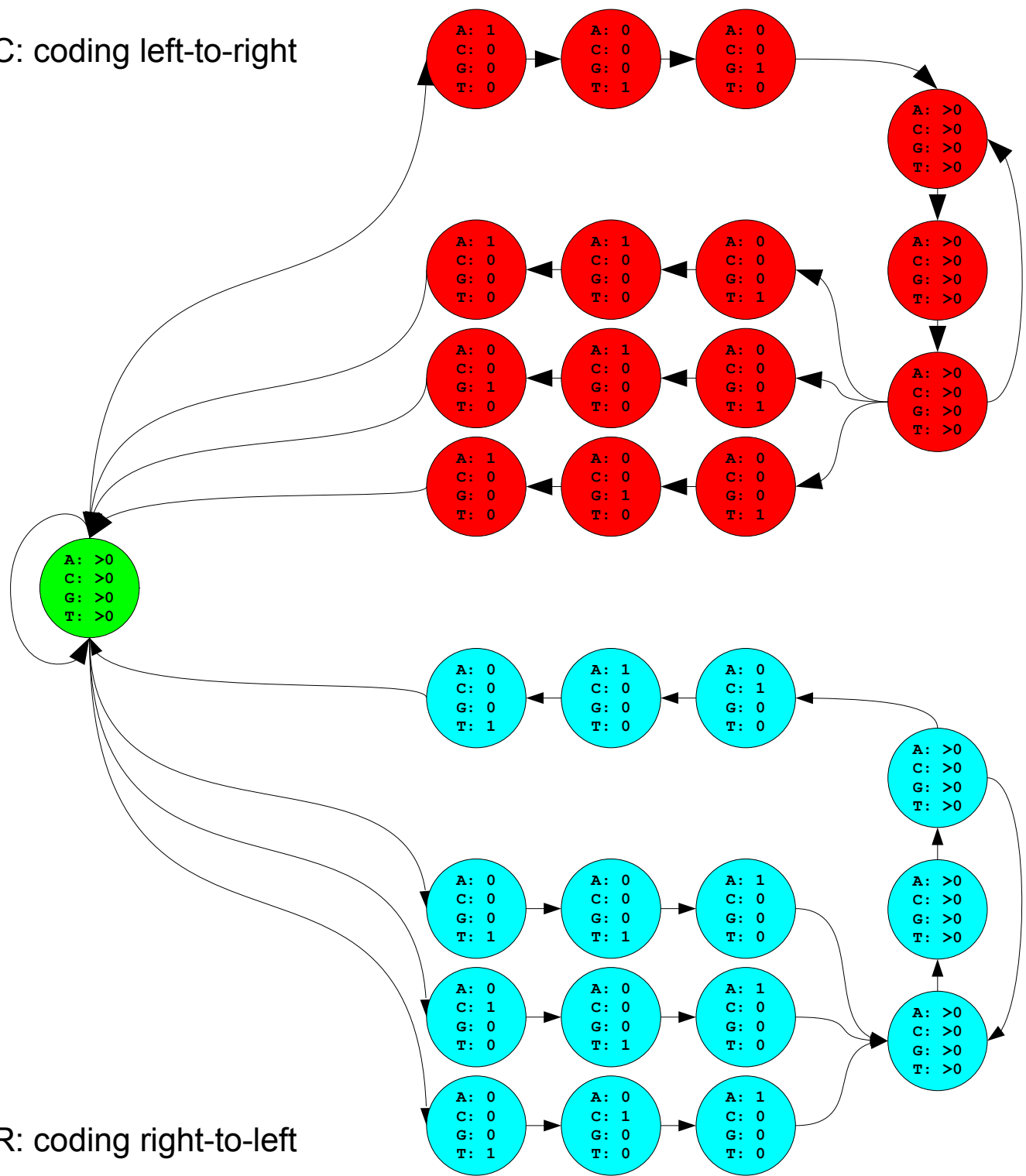
A "codon" model



C: coding left-to-right

N: Non-coding

R: coding right-to-left



ACs between predictions and true annotations sorted by average

	Genome 6	Genome 7	Genome 8	Genome 9	Genome 10	Average
3	0,6313	0,6378	0,6007	0,5310	0,4762	0,5754
22	0,6313	0,6378	0,6007	0,5310	0,4762	0,5754
9	0,6053	0,6147	0,5708	0,4919	0,4197	0,5405
10	0,4628	0,4587	0,4245	0,3244	0,3611	0,4065
21	0,4619	0,4582	0,4219	0,3239	0,3603	0,4052
13	0,4539	0,4547	0,4159	0,3121	0,3650	0,4003
11	0,4522	0,4504	0,4172	0,3139	0,3648	0,3997
1	0,4520	0,4529	0,4165	0,3122	0,3646	0,3996
25	0,4543	0,4515	0,4144	0,3133	0,3596	0,3986
18	0,4521	0,4481	0,4112	0,3095	0,3619	0,3966
15	0,4501	0,4449	0,4187	0,2974	0,3586	0,3939
24	0,4237	0,4326	0,4063	0,2769	0,3816	0,3842
14	0,4359	0,4383	0,3974	0,2642	0,3467	0,3765
2	0,4286	0,4548	0,3725	0,2473	0,3733	0,3753
19	0,3856	0,4075	0,3727	0,2926	0,3059	0,3529
23	0,3854	0,4073	0,3551	0,2336	0,3710	0,3505
5	0,3850	0,4008	0,3619	0,2559	0,3119	0,3431
6	0,3821	0,3914	0,3504	0,2365	0,3263	0,3373
16	0,2917	0,3462	0,2618	0,2240	0,3213	0,2890
7	0,0247	0,0978	0,0311	0,0642	0,1441	0,0724
4	0,0254	0,0939	0,0253	0,0631	0,1387	0,0693
20	0,0775	0,0526	0,1344	-0,0246	0,0765	0,0633
12	0,0344	0,0820	0,0089	0,0549	0,1144	0,0589
17	-0,2748	-0,2571	-0,2635	-0,2657	-0,3075	-0,2737
8						

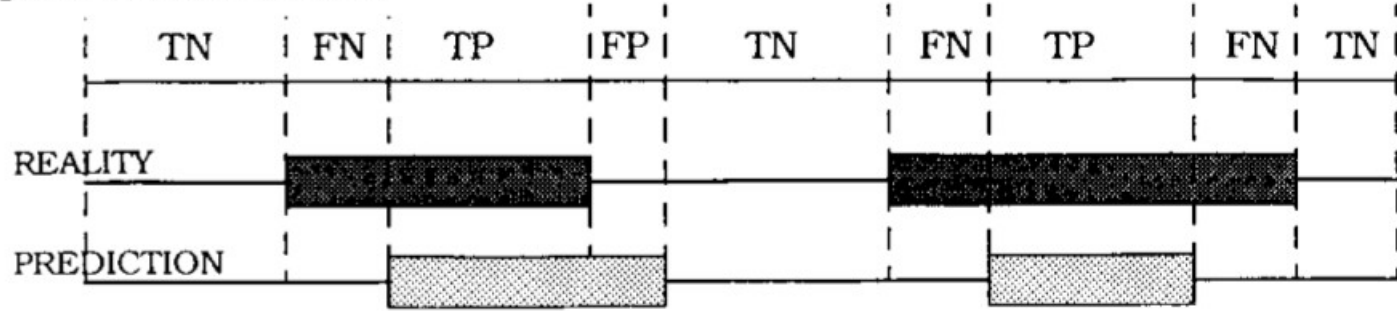
Codon-models

Models inspired by the model presented in class maybe allowing more start/stop codons.

Problems with training or prediction

Evaluating performance

Nucleotide Level



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

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

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

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

History and applications of HMMs

History of HMMs

Hidden Markov Models were introduced in statistical papers by Leonard E. Baum and others in the late 1960s. One of the first applications of HMMs was speech recognition in the mid-1970s.

In the late 1980s, HMMs were applied to the analysis of biological sequences. Since then, many applications in bioinformatics...

Applications of HMMs in bioinformatics

prediction of protein-coding regions in genome sequences

modeling families of related DNA or protein sequences

prediction of secondary structure elements in proteins

... and many others ...