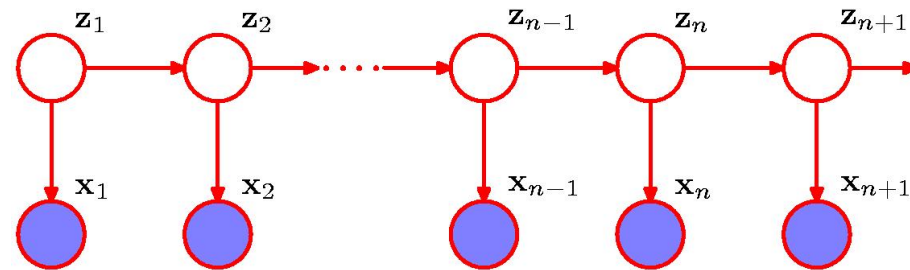


Hidden Markov Models

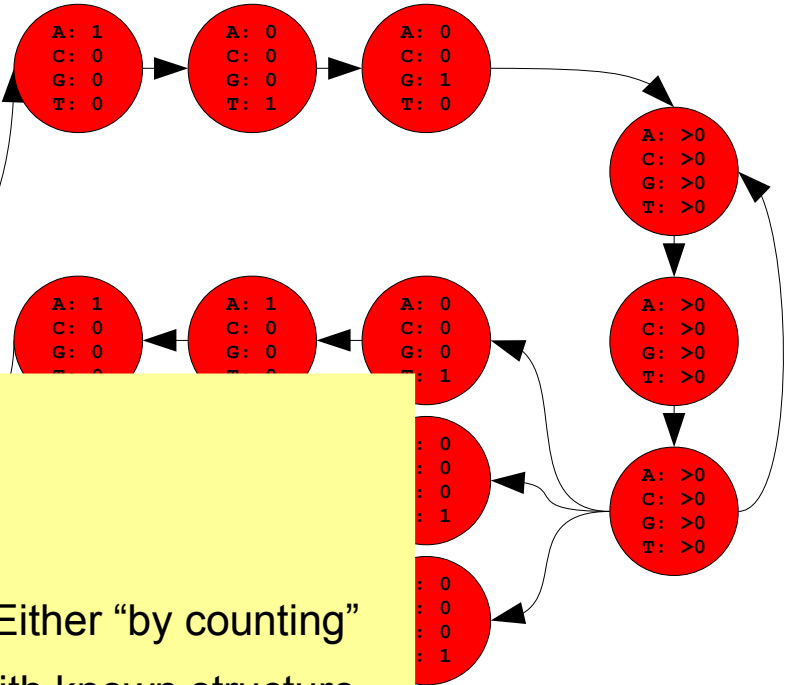
Hand In 3 - Gene finding



Even more biology

There can be genes in both directions

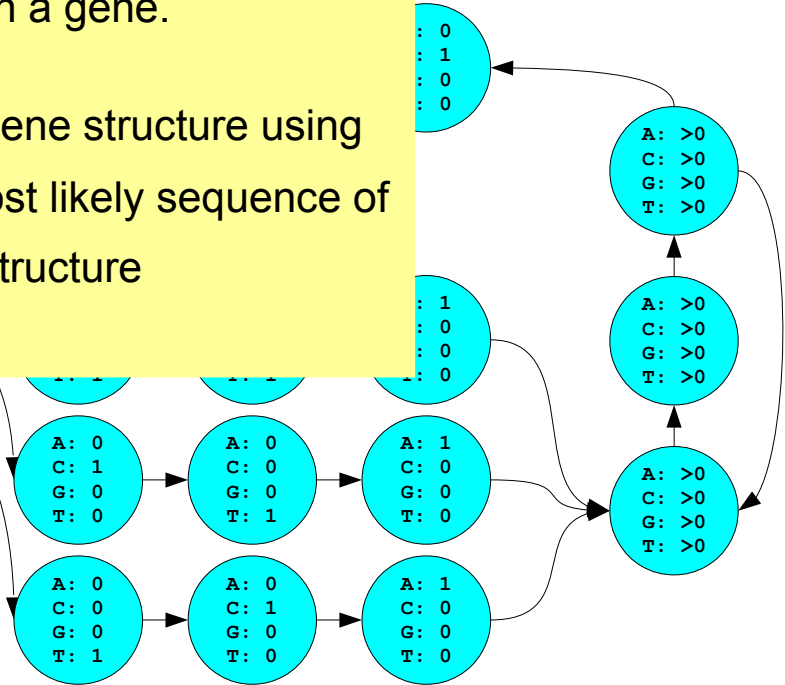
C: coding left-to-right



Gene finding

- Select initial model structure
- 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



R: coding right-to-left


$$\pi_N = 1$$

$$\pi_C = 0$$

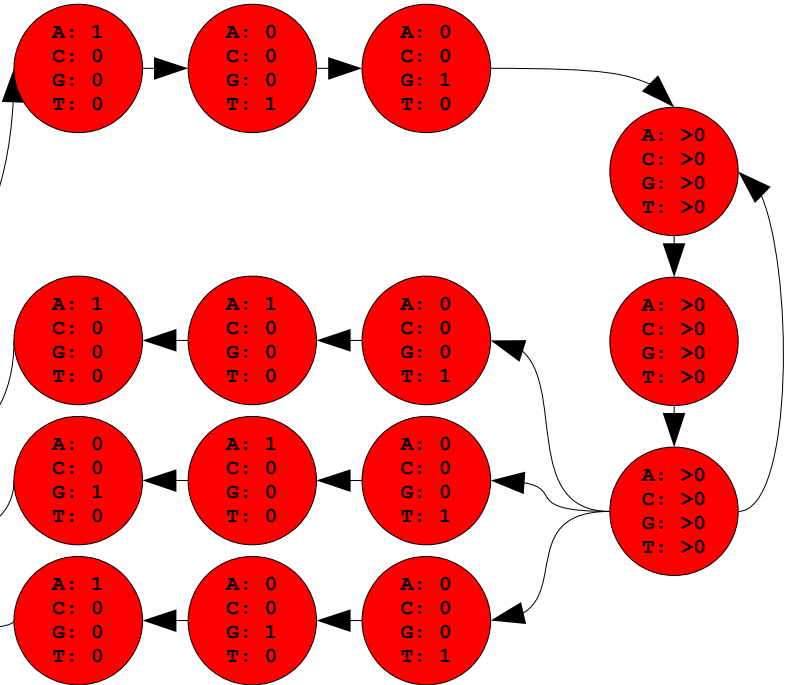
Selecting an initial model

Even more biology

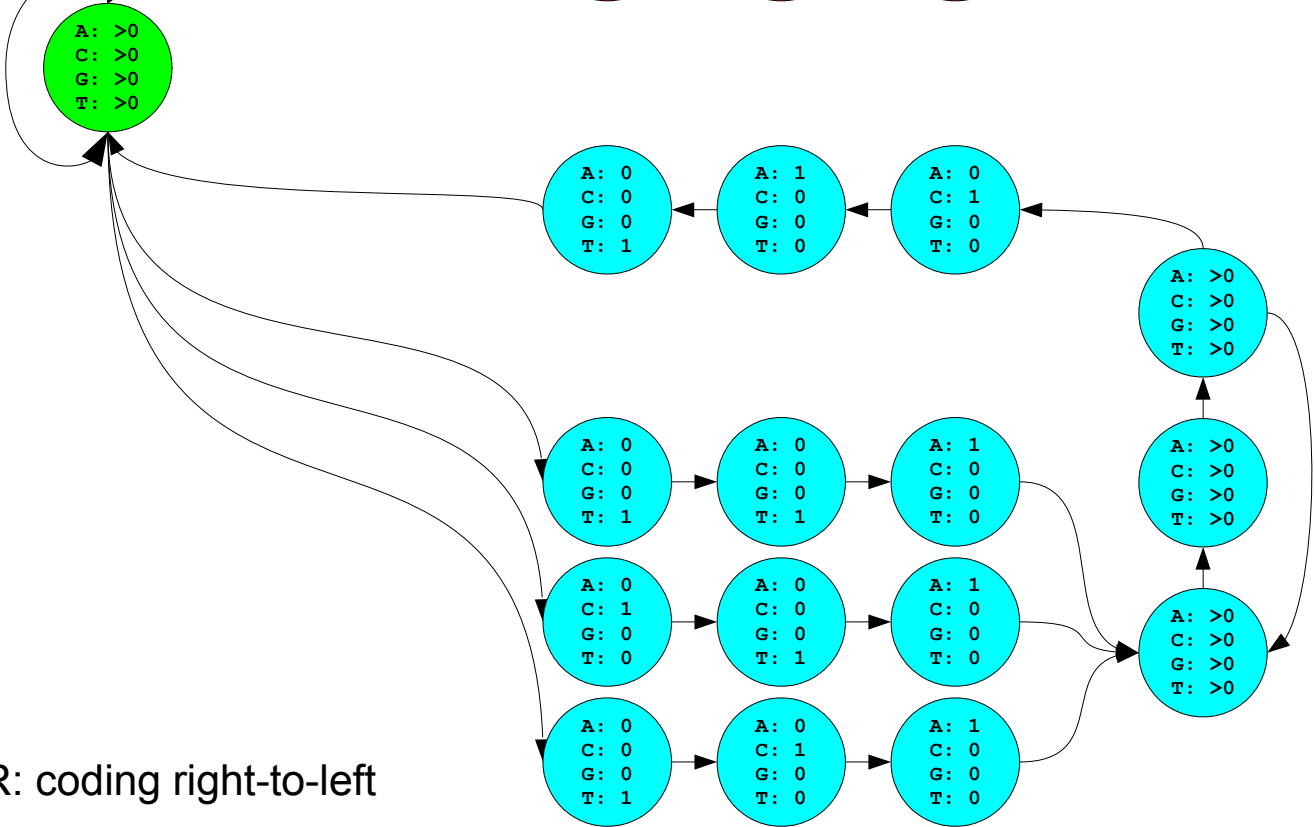
There can be genes in both directions



C: coding left-to-right



N: Non-coding



R: coding right-to-left

$$\pi_N = 1$$

$$\pi_C = 0$$

Analysis of the training data

```
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)
```

We observe

3 typical start-codons

3 stop-codons

We may ignore rare start and stop codons, i.e. a gene that starts (or ends) with a ignored start (or stop) codon does not correspond to a path in your model.

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']
```

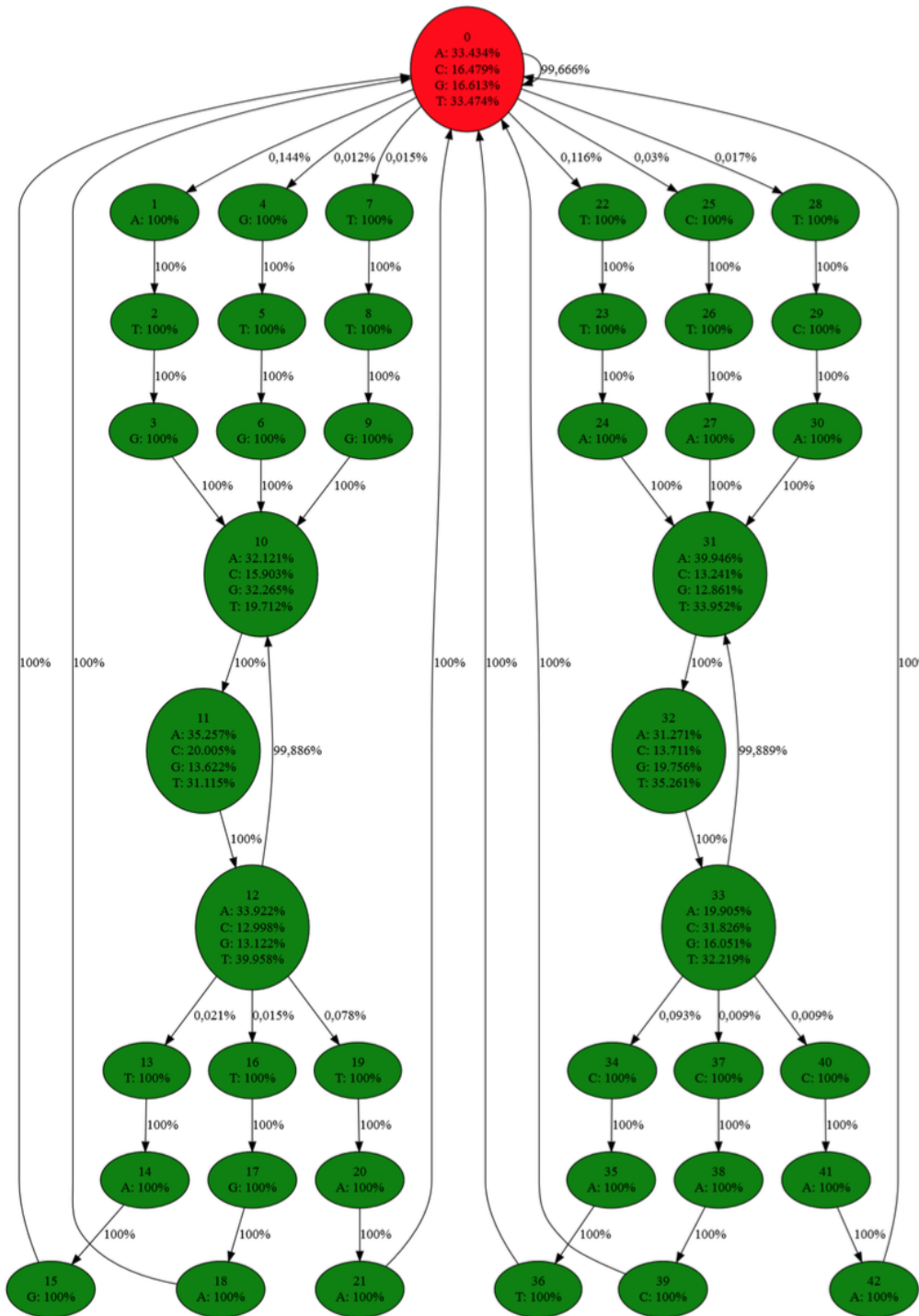
Reversed stop-codon in reversed genes:

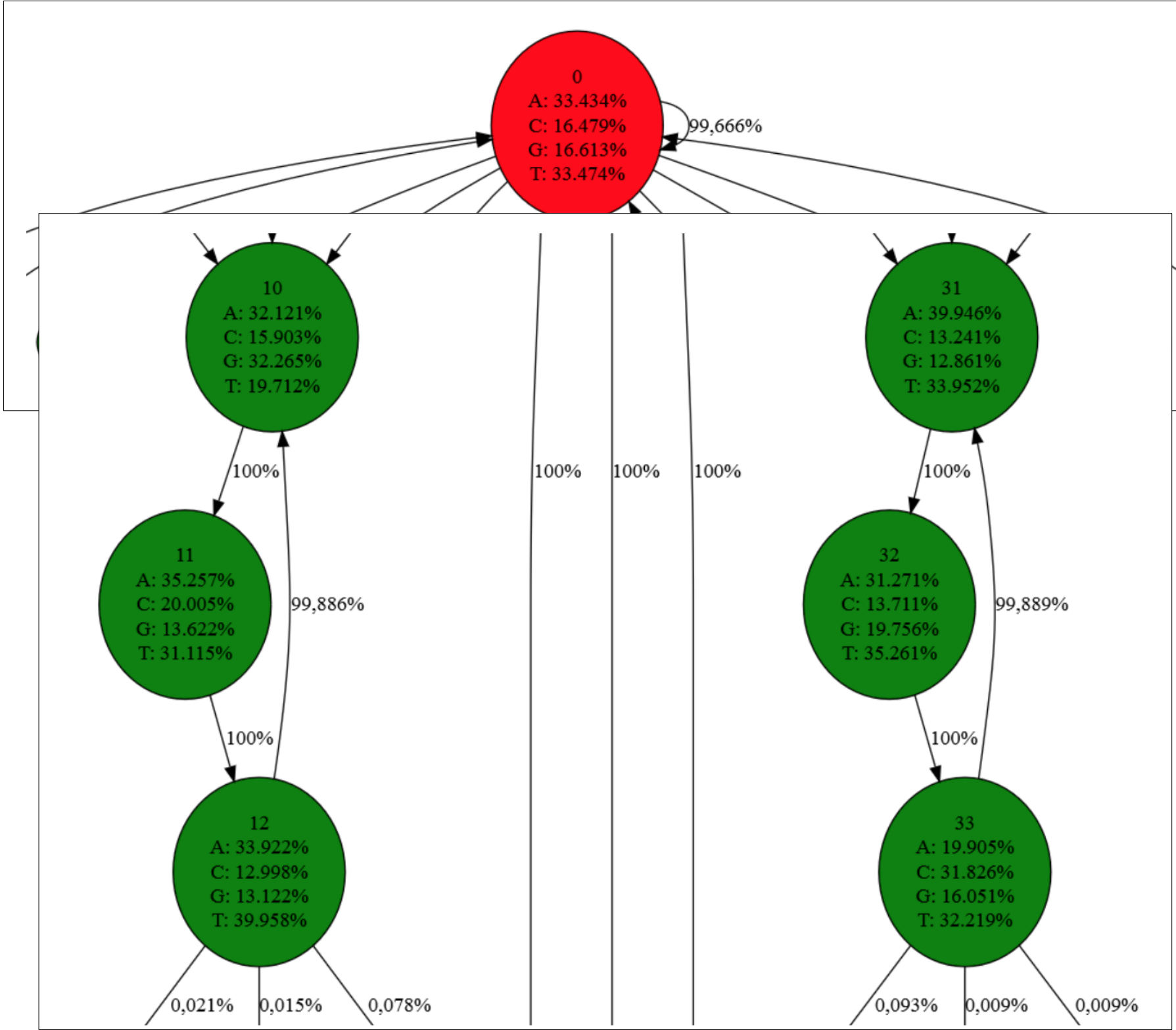
```
TTA (reverse-complement: TAA) [6596, 'NRRR']
CTA (reverse-complement: TAG) [2014, 'NRRR']
TCA (reverse-complement: TGA) [1148, 'NRRR']
```

Reversed start-codon in reversed genes:

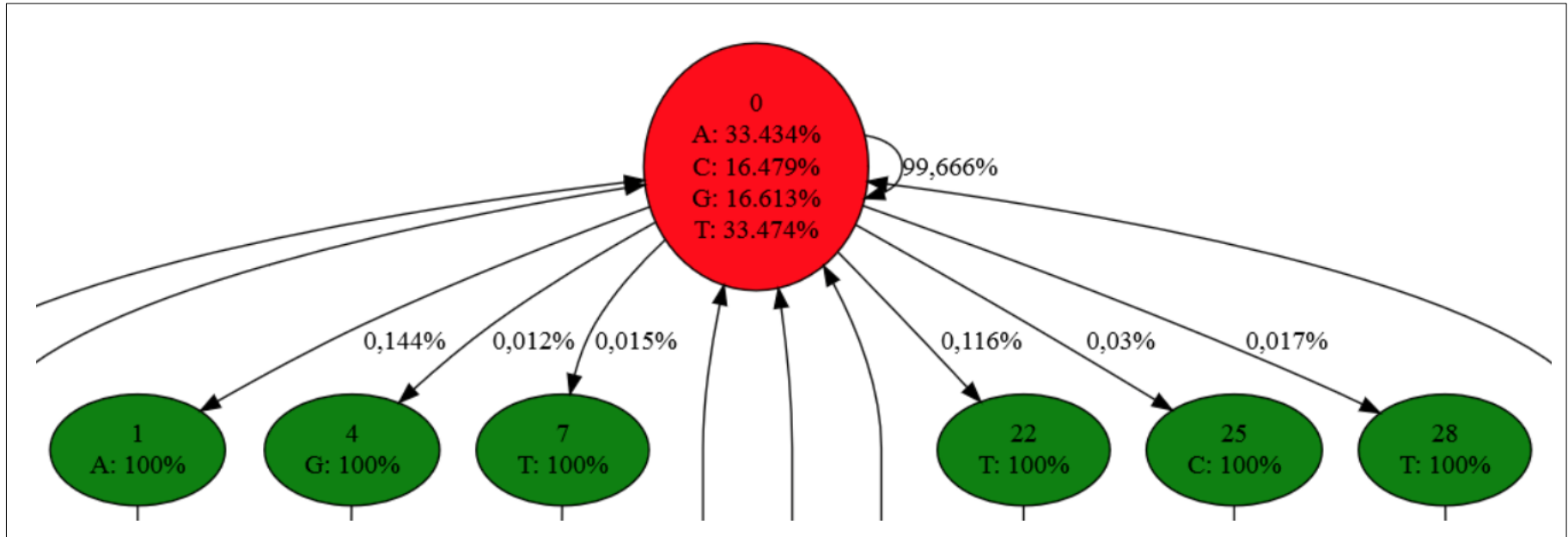
```
TAT (reverse-complement: ATA) [2, 'RRRN']
ATG (reverse-complement: CAT) [1, 'RRRN']
GAT (reverse-complement: ATC) [1, 'RRRN']
CAT (reverse-complement: ATG) [8077, 'RRRN']
AAT (reverse-complement: ATT) [4, 'RRRN']
TAC (reverse-complement: GTA) [1, 'RRRN']
CAC (reverse-complement: GTG) [715, 'RRRN']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

A typical model





Training by counting



In theory, we are given (\mathbf{X}, \mathbf{Z}) pairs, but we are given $(\mathbf{X}, \text{"Z"})$ pairs, where the NCR-annotations have to be translated to \mathbf{Z} s.

Also, we may ignore rare start and stop codons, i.e. a gene that starts (or ends) with a ignored start (or stop) codon does not correspond to a path in your model.

Training by counting – Typical solution

To set the transition probabilities:

$N \rightarrow N$	$N \rightarrow CCC$, where CCC is ATG	$N \rightarrow RRR$, where RRR is TTA
	$N \rightarrow CCC$, where CCC is GTG	$N \rightarrow RRR$, where RRR is CTA
	$N \rightarrow CCC$, where CCC is TTG	$N \rightarrow RRR$, where RRR is TCA

We count:

$\#(N \rightarrow N)$ = no. of occurrences of “NN” our annotations.

$\#(N \rightarrow CCC, \text{ where } CCC \text{ is } XYZ)$ = no. of occurrence “NCCC” in our annotations, where CCC is an annotation of XYZ (in our training data).

$\#(N \rightarrow RRR, \text{ where } RRR \text{ is } XYZ)$ = no. of occurrence “NRRR” in our annotations where, RRR is an annotation of XYZ (in our training data).

We compute:

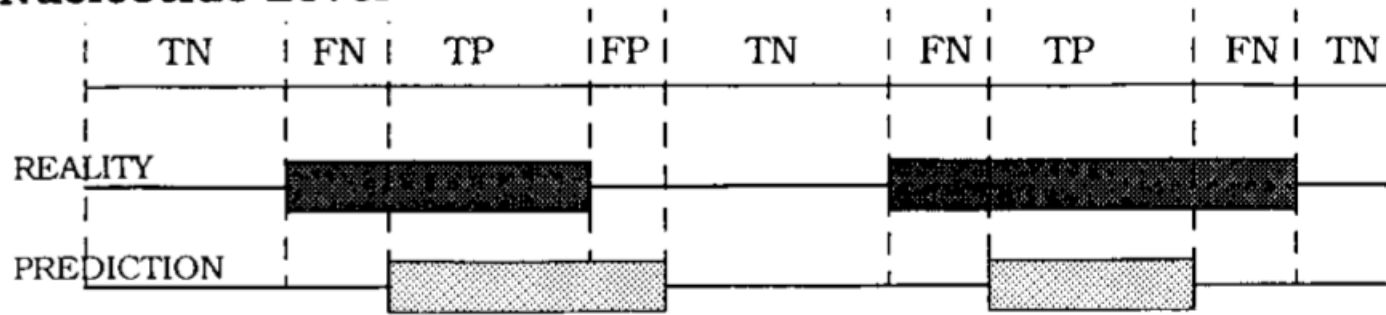
$\text{Total} = \#(N \rightarrow N) + \#(N \rightarrow CCC \text{ where } CCC \text{ is } XYZ) + \#(N \rightarrow RRR, \text{ where } RRR \text{ is } XYZ)$

We set:

$P(N \rightarrow X) = \#(N \rightarrow X) / \text{Total}$ for each of the 7 transitions

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

compare_anns.py

```
$ python compare_anns.py true-ann6.fa pred-ann6.fa
```

```
Cs (tp=757332, fp=164766, tn=305197, fn=57217): Sn = 0.9298, Sp = 0.8213, AC = 0.6213
```

```
Rs (tp=715865, fp=127462, tn=304830, fn=57584): Sn = 0.9255, Sp = 0.8489, AC = 0.6603
```

```
Both (tp=1473197, fp=292228, tn=247613, fn=114801): Sn = 0.9277, Sp = 0.8345, AC = 0.4520
```

k-fold cross validation

For each genome 1 to 5, you train your model by training-by-counting on the remaining 4 genomes, and predict the gene structure on the genome you picked. You compute and report the approximate correlation coefficient (AC) between your prediction and the true annotation using the small python program `compare_anns.py`. Include a table with the computed ACs in your report.

	Round 1	Round 2	Round 3	Round 4	Round 5
Genome 1	Validate	Train	Train	Train	Train
Genome 2	Train	Validate	Train	Train	Train
Genome 3	Train	Train	Validate	Train	Train
Genome 4	Train	Train	Train	Validate	Train
Genome 5	Train	Train	Train	Train	Validate

Table 1: ACs between the predictions on the genomes 1 to 5 and their true annotation. The gene predictor was trained by the remaining four genomes respectively.

	AC		
	Only Cs	Only Rs	Both
Genome 1	0.5983	0.6470	0.4347
Genome 2	0.6305	0.6529	0.4510
Genome 3	0.6552	0.6510	0.4805
Genome 4	0.6240	0.6002	0.4079
Genome 5	0.6550	0.6027	0.4302

	Round 1	Round 2	Round 3	Round 4	Round 5
Genome 1	Validate	Train	Train	Train	Train
Genome 2	Train	Validate	Train	Train	Train
Genome 3	Train	Train	Validate	Train	Train
Genome 4	Train	Train	Train	Validate	Train
Genome 5	Train	Train	Train	Train	Validate

Performance on genome 6-10

The screenshot shows a web browser window with the URL `services.birc.au.dk/genefinder-verifier/`. The page title is "GeneFinder Verifier". Below the title, there is a "Welcome!" section with a paragraph explaining the tool's purpose: comparing predicted gene structures of genomes 6 to 10 against true gene structures. A code block shows the expected input format for a FASTA file, listing prediction files for genomes 6 through 10. Below this, a section titled "The output will look like this:" displays performance metrics for each genome (6-10). For each genome, three lines of data are shown for "Cs", "Rs", and "Both" methods, each with values for Sensitivity (Sn), Specificity (Sp), and Accuracy (AC). At the bottom, there is a blue upload area with the text "Upload your sequences here." and a note about file size. A file selection button labeled "Choose File" is currently inactive, and an "Upload" button is visible below it.

GeneFinder Verifier

Welcome!

Here you can compare the predicted genes structures of genome 6 to 10 against the true gene structures. The predicted gene structures are given in a single fasta file that must have this format (you can leave out the entries for the genomes that you do not want to be part of the comparison):

```
> pred-ann6
[PREDICTION OF GENOME6]
> pred-ann7
[PREDICTION OF GENOME7]
> pred-ann8
[PREDICTION OF GENOME8]
> pred-ann9
[PREDICTION OF GENOME9]
> pred-ann10
[PREDICTION OF GENOME10]
```

The output will look like this:

```
Genome 6
Cs (tp=757332, fp=164766, tn=305197, fn=57217): Sn = 0.9298, Sp = 0.8213, AC = 0.6213
Rs (tp=715865, fp=127462, tn=304830, fn=57584): Sn = 0.9255, Sp = 0.8489, AC = 0.6603
Both (tp=1473197, fp=292228, tn=247613, fn=114801): Sn = 0.9277, Sp = 0.8345, AC = 0.4520
Genome 7
Cs (tp=868820, fp=236008, tn=517048, fn=79049): Sn = 0.9166, Sp = 0.7864, AC = 0.6285
Rs (tp=815026, fp=226580, tn=511963, fn=84134): Sn = 0.9064, Sp = 0.7825, AC = 0.6205
Both (tp=1683846, fp=462588, tn=432914, fn=163183): Sn = 0.9117, Sp = 0.7845, AC = 0.4529
Genome 8
Cs (tp=705403, fp=137180, tn=351159, fn=74782): Sn = 0.9041, Sp = 0.8372, AC = 0.6424
Rs (tp=607762, fp=169829, tn=351738, fn=74203): Sn = 0.8912, Sp = 0.7816, AC = 0.5865
Both (tp=1313165, fp=307009, tn=276956, fn=148985): Sn = 0.8981, Sp = 0.8105, AC = 0.4166
Genome 9
Cs (tp=776640, fp=203664, tn=340882, fn=88415): Sn = 0.8978, Sp = 0.7922, AC = 0.5550
Rs (tp=759048, fp=219786, tn=336181, fn=93116): Sn = 0.8907, Sp = 0.7755, AC = 0.5270
Both (tp=1535688, fp=423450, tn=247766, fn=181531): Sn = 0.8943, Sp = 0.7839, AC = 0.3122
Genome 10
Cs (tp=612457, fp=106124, tn=253878, fn=88014): Sn = 0.8744, Sp = 0.8523, AC = 0.5872
Rs (tp=371869, fp=138143, tn=291605, fn=50287): Sn = 0.8809, Sp = 0.7291, AC = 0.5707
Both (tp=984326, fp=244267, tn=203591, fn=138301): Sn = 0.8768, Sp = 0.8012, AC = 0.3640
```

Upload your sequences here.

Note that the size of a fasta file containing the predictions of all five genomes is about 10 Mb, so uploading and comparison takes some seconds.

File: no file selected

<https://services.birc.au.dk/genefinder-verifier/>

Performance on genome 6-10

Input is a fasta-file containing your predictions. It is about 10 Mb big!

```
> pred-ann6  
[ PREDICTION OF GENOME6 ]  
> pred-ann7  
[ PREDICTION OF GENOME7 ]  
> pred-ann8  
[ PREDICTION OF GENOME8 ]  
> pred-ann9  
[ PREDICTION OF GENOME9 ]  
> pred-ann10  
[ PREDICTION OF GENOME10 ]
```

```
$ cat pred-ann6.fa pred-ann7.fa pred-ann8.fa pred-ann9.fa pred-ann10.fa > pred-ann6-10.fa
```

Beware of the newlines!

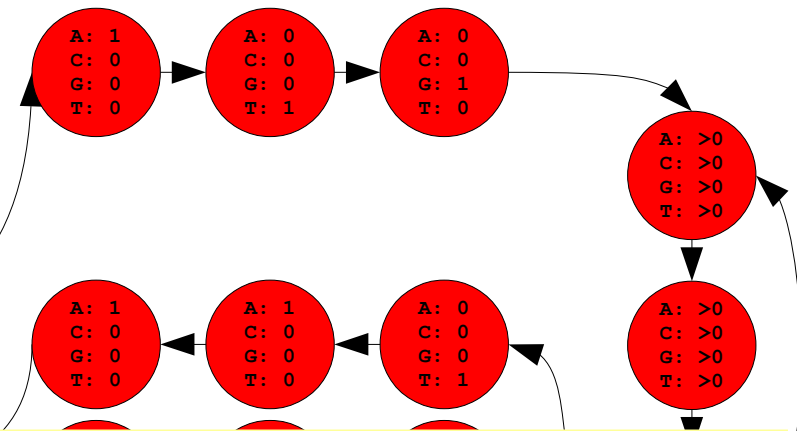
pred-annX.fa should end with a newline in order for the above to work

Last year

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	Codon-models
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	0,4063	Models inspired by the model presented in class maybe allowing more start/stop codons.
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	0,3643	0,3997	
1	0,4520	0,4529	0,4165	0,3122	0,3645	0,3996	
25	0,4543	0,4515	0,4144	0,3133	0,3595	0,3986	
18	0,4521	0,4481	0,4112	0,3095	0,3619	0,3966	
15	0,4501	0,4449	0,4187	0,2974	0,3585	0,3939	
24	0,4237	0,4326	0,4063	0,2769	0,3815	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	Problems with training or prediction
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							

C: coding left-to-right

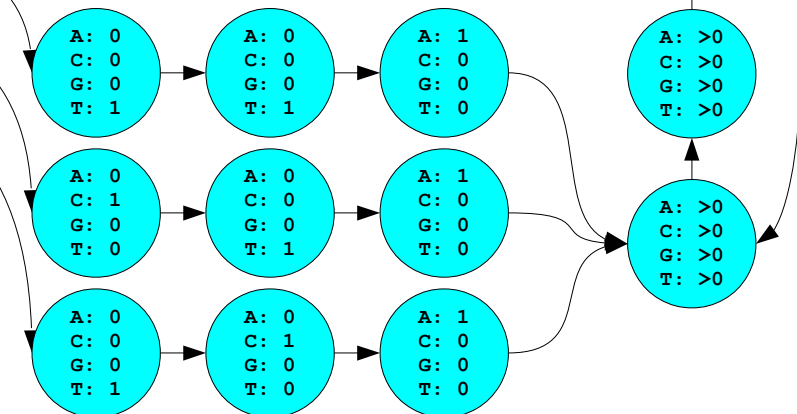


A “standard” model

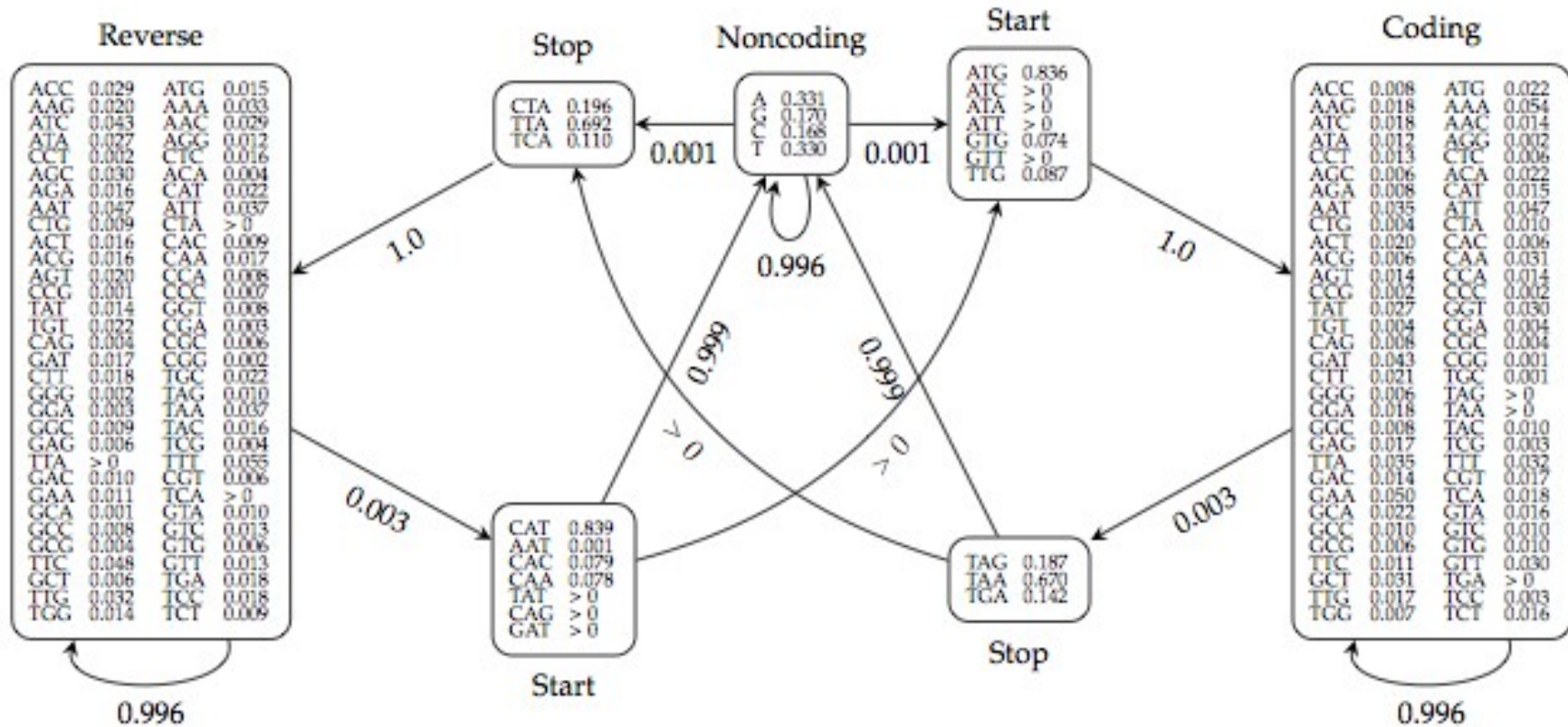
What makes it possible to distinguish between coding and non-coding regions is the difference in nucleotide frequencies.

The explicit model of start- and stop-codons makes sure that we only make “syntactically” correct transitions between coding and non-coding regions.

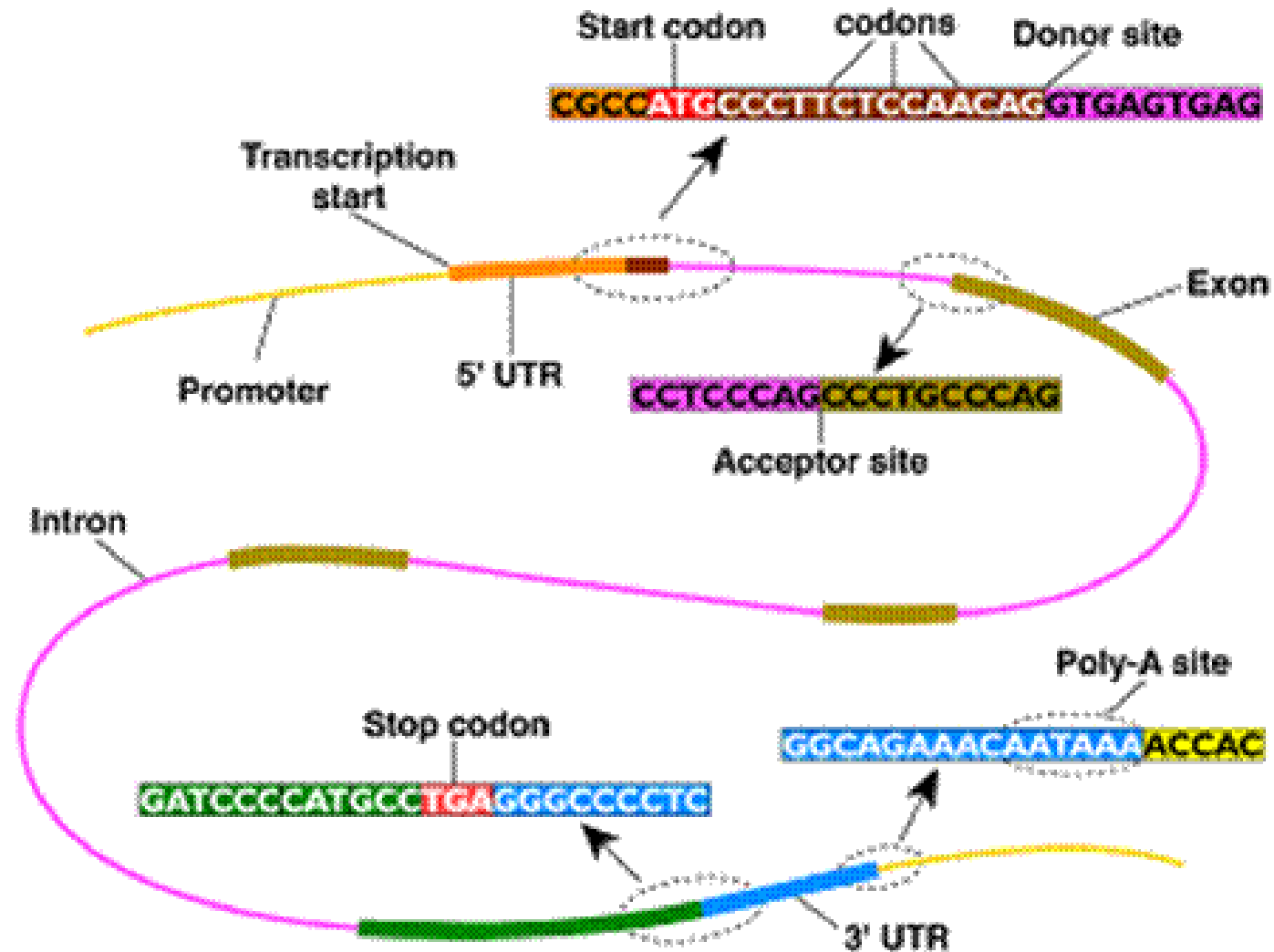
R: coding right-to-left



A "codon" model

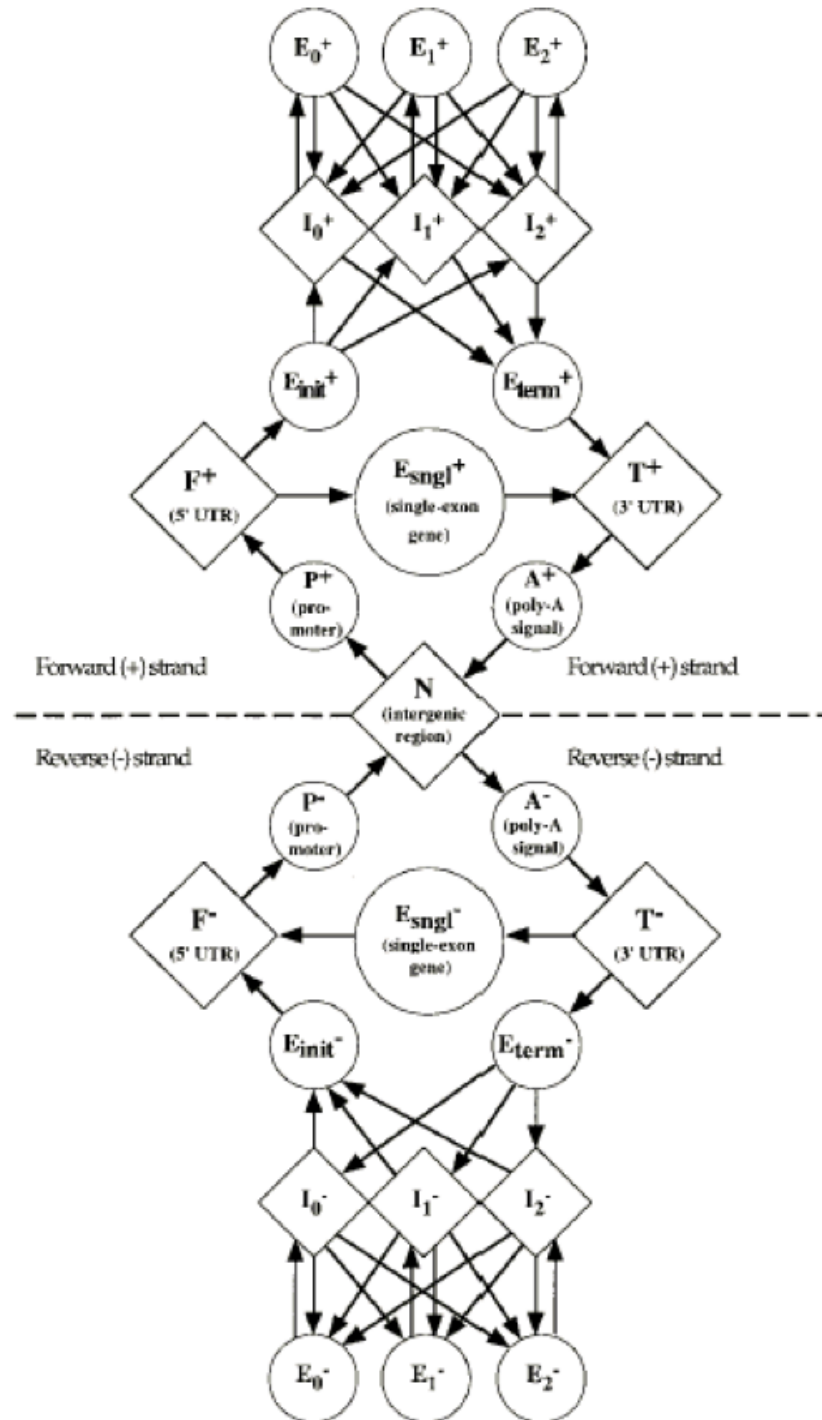


More complex gene finding problems



Eukaryotic gene structure

Eukaryotic gene structure in both directions (GenScan)



Exam

Overall: Oral, 20 min all included, 10 min for your **prepared** presentation, 5 min for discussion/question, you may be interrupted along the way. See Blackboard for all the details.

HMM Question: Hidden Markov Models (Basic algorithms and applications, Building models and selecting model parameters)

What is a HMM, explain/show parameters and assumptions

How to compute the joint probability $P(X,Z)$

Basic problems: Decoding and selecting models parameter, e,g:

How to compute a Viterbi decoding, Z^* , where $P(Z^*) = \max_z P(X,Z)$

How to compute $P(X) = \sum_z P(X,Z)$ using the forward algorithm

How to compute a posterior decoding using the forward and backward algorithms

How to select parameters from (X,Z) using training-by-counting, or from X using EM

Remember to be **precise**, and remember **what, why, and how**, when you explain