

# PROBABILISTIC MODELS: HIDDEN MARKOV MODELS



**PROLOGUE:**

**Pitfalls of standard alignments**

# Scoring a pairwise alignment

A: ALA**E**VLIRLIT**K**LYP  
 B: ASA**K**HLNRLIT**E**LYP

$$Score(A, B) = \sum s(A^i, B^i)$$

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	-4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2	
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

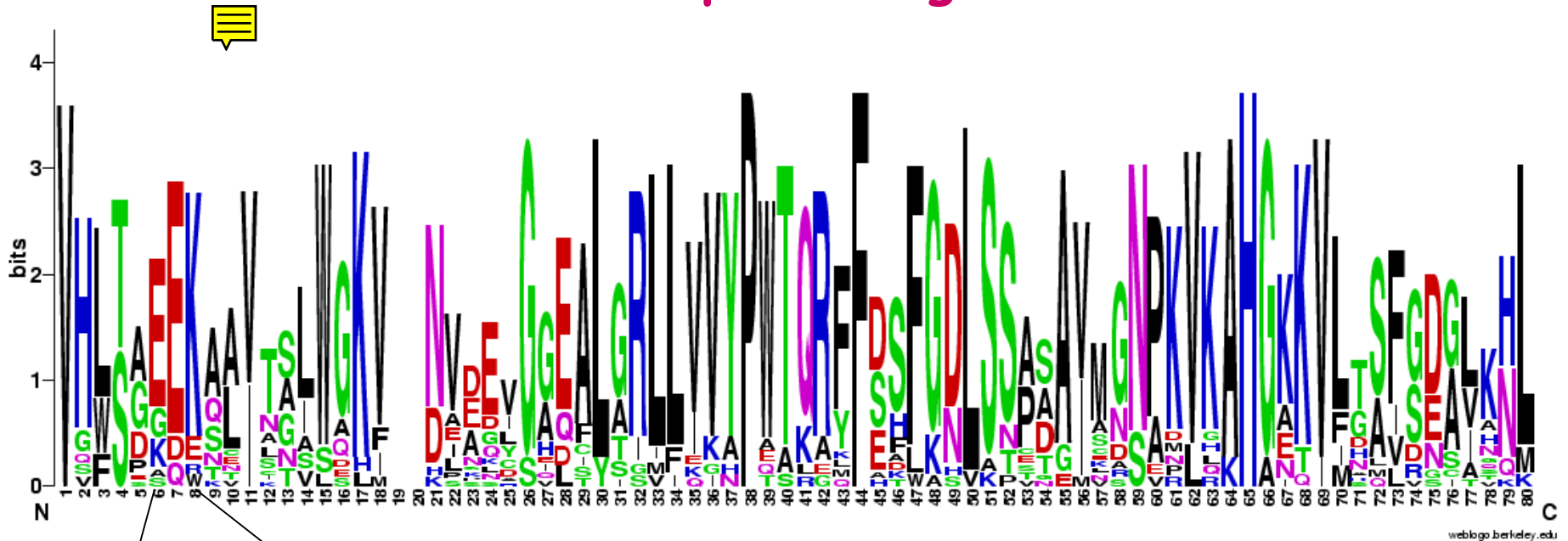
Blosum62

# Alignment of a family (globins)

		10	20	30	40	50	60	70	80													
lqb1	pea/1-471	-GFTDKQ	EALVNSSSE	-FKQNL	PGYSILFY	TIVLEK	APAAKGL	SFLKD---	TAGVED	SPKLQAHAEQ	VFGLVRD	SAAQL	.									
lqb1	vicfa/1-471	-GFTEKQ	EALVNSSSQ	LFKQNP	SNYSVLFY	TIILQK	APTAKAMF	SFLKD---	SAGVVVD	SPKLGAHA	EKVFGM	VRD	SAVQL	.								
hbb	speci/1-471	VHLS	DGEKNA	ISTAWG	KV--HAA	EVGA	EALGRLL	VVYPWT	QRF	FDSFGDL	SSASAV	MGNAKV	KAHGKK	VID	SFSNGL	KHL	.					
hbb	speto/1-471	VHLT	DGEKNA	ISTAWG	KV--NAA	EIGA	EALGRLL	VVYPWT	QRF	FDSFGDL	SSASAV	MGNAKV	KAHGKK	VID	SFSNGL	KHL	.					
hbb	equhe/1-471	VQLS	GEEKAA	AVLAL	WDKV--	NEEE	VGGEAL	GRLLV	VYPWT	QRF	FDSFGDL	SNPAA	VMGNP	KVKA	HGKKV	LH	SFGE	GVHHL	.			
hbb	sunmu/1-471	VHLS	GEEKAC	VTGLW	GKV--	NEDE	VGA	EALGRLL	VVYPWT	QRF	FDSFGDL	SSASAV	MGNP	KVKA	HGKKV	LH	SLGEG	VANL	.			
hbb	tupql/1-471	VHLS	GEEKAA	AVTGL	WGKV--	DLEK	VGGQ	SLGSL	LIVYPWT	QRF	FDSFGDL	SSPSAV	MSNP	KVKA	HGKKV	L	TSFSD	G	L	NHL	.	
hbb	calar/1-471	VHLT	GEEKS	AVTAL	WGKV--	NVDE	VGGEAL	GRLLV	VYPWT	QRF	FESFGDL	STPDA	VMNNP	KVKA	HGKKV	L	GAFSD	G	L	THL	.	
hbb	mansp/1-471	VHLT	PEEKT	AVTTL	WGKV--	NVDE	VGGEAL	GRLLV	VYPWT	QRF	FDSFGDL	SSPDA	VMGNP	KVKA	HGKKV	L	GAFSD	G	L	NHL	.	
hbb	rabit/1-471	VHLS	SSEKS	AVTAL	WGKV--	NVEE	VGGEAL	GRLLV	VYPWT	QRF	FESFGDL	SSANAV	MMNP	KVKA	HGKKV	L	A	AFSEGL	S	HL	.	
hbb	ursma/1-471	VHLT	GEEKS	LVTL	WGKV--	NVDE	VGGEAL	GRLLV	VYPWT	QRF	FDSFGDL	SSADA	IMNP	KVKA	HGKKV	L	NSFSD	G	L	KNL	.	
hbb	triin/1-471	VHLT	PEEK	ALVIG	LWAKV--	NVKE	YGGEAL	GRLLV	VYPWT	QRF	FEHFGDL	SSASA	IMNP	KVKA	HGEKV	FTS	FGD	G	L	KHL	.	
hbb	ornan/1-471	VHLS	GGEKS	AVTNL	WGKV--	NINEL	GGEAL	GRLLV	VYPWT	QRF	FEAFGDL	SSAGAV	MGNP	KVKA	HGAKV	L	TSFGD	A	L	KNL	.	
hbb	tacac/1-471	VHLS	GSEKT	AVTNL	WGHV--	NVNE	LGG	EALGRLL	VYPWT	QRF	FESFGDL	SSADA	VMGN	AKVKA	HGAKV	L	TSFGD	A	L	KNL	.	
hbe	ponpy/1-471	VHFT	AEEKA	AAVTS	LWSKM--	NVEE	AGGEAL	GRLLV	VYPWT	QRF	FDSFGNL	SSPSA	ILGNP	KVKA	HGKKV	L	TSFGD	A	I	KNM	.	
hbb	colli/1-471	VHWS	AEEK	QLITS	IWGKV--	NVAD	CGAEAL	ARLLI	VYPWT	QRF	FSSFGNL	SSATA	ISGNP	NVKA	HGKKV	L	TSFGD	A	V	KNL	.	
hbb	larri/1-471	VHWS	AEEK	QLIT	GLWGKV--	NVAD	CGAEAL	ARLLI	VYPWT	QRF	FASFGNL	SSPTA	INGN	PMVRA	HGKKV	L	TSFGE	A	V	KNL	.	
hbb1	varex/1-471	VHWT	AEEK	QLIC	SLWGKI--	DVGL	IGGET	LAGLL	VIYPWT	QRF	FSHF	GNLSS	PTA	IAGN	PRVKA	HGKKV	L	TSFGD	A	I	KNL	.
hbb2	xentr/1-471	VHWT	AEEK	KATIA	SVWGKV--	DIEQ	DGHD	ALSRL	LVYPWT	QRY	FSSFGNL	SNVSA	VSGN	VKVKA	HGNKV	L	SAVGS	A	I	QHL	.	
hbb1	ranca/1-471	VHWT	AEEK	KAVIN	SVWQKV--	DVEQ	DGHEAL	TRLF	IVYPWT	QRY	FSTFGDL	SSPAA	IAGN	PKVHA	HGKKI	L	G	AIDNA	I	HNL	.	
hbb2	tracr/1-471	VHLT	AEDR	KEIAA	ILGKV--	NVDS	LGGQ	CLARL	IVNPN	WRRY	FHDF	GDLSS	CD	AICRN	PKVLA	HGAKV	MRSI	VEAT	K	HL	.	
hba4	salir/1-471	-SL	SAKDK	ANVKA	IWGKIL	PKSDE	IGEQA	LSRML	VVYPQ	TKAYF	SHMAS	VAP----	GSAP	VKKHG	ITIM	NQID	DCVGHM	.	.	.	.	
myg	escqi/1-471	-VL	SDAE	WQLVL	NIWAK	VEADV	AGHGQD	ILIR	LFKG	HPET	LEKFD	KFKHL	KTEA	EMKASE	DLKKH	GNTVL	TALG	GIL	KKK	.	.	

Different positions are not equivalent

# Sequence logos



<http://weblogo.berkeley.edu/cache/file5h2DWc.png>



The substitution score IN A FAMILY should depend on the position (the same for gaps)

For modelling families we need more flexible tools

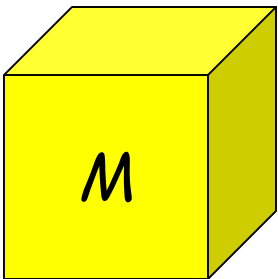
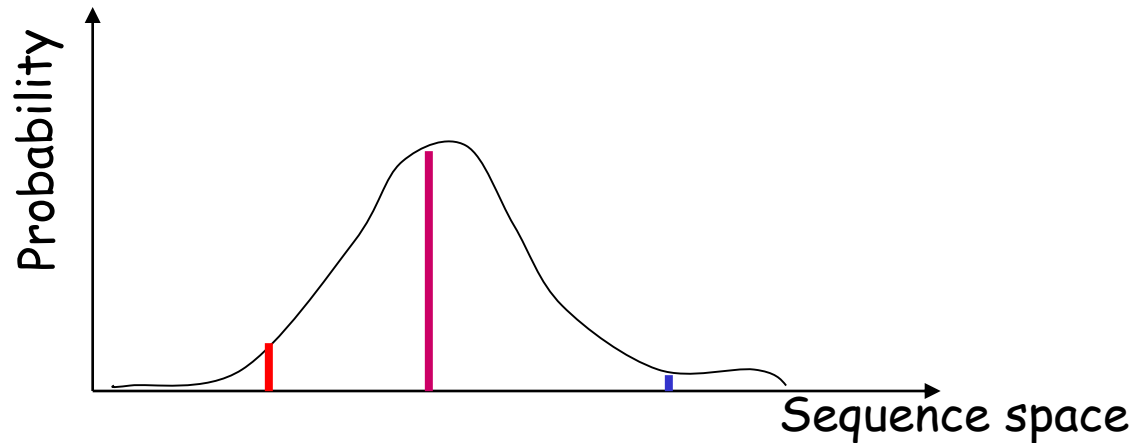
# Probabilistic Models for Biological Sequences

- What are they?

# Probabilistic models for sequences

## Generative definition:

- Objects producing different outcomes (sequences) with different probabilities
- The probability distribution over the sequences space determines the model specificity



*Generates  $s_i$  with probability  $P(s_i | M)$*   
e.g.:  $M$  is the representation of the family of globins

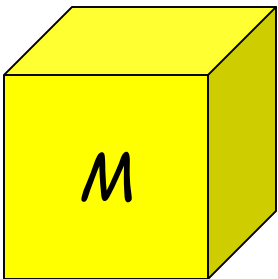
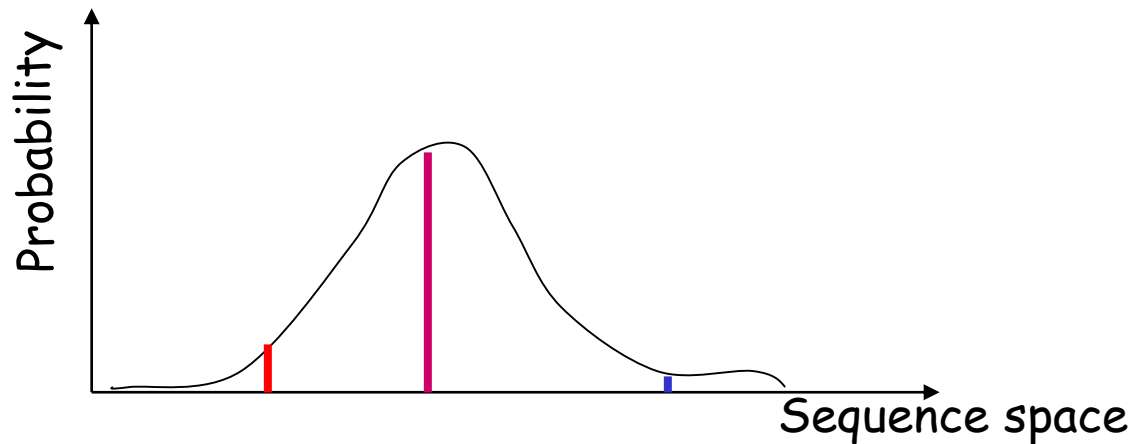
# Probabilistic models for sequences

*We don't need a generator of new biological sequences*

the generative definition is useful as operative definition

*Associative definition:*

- Objects that, given an outcome (sequence), compute a probability value



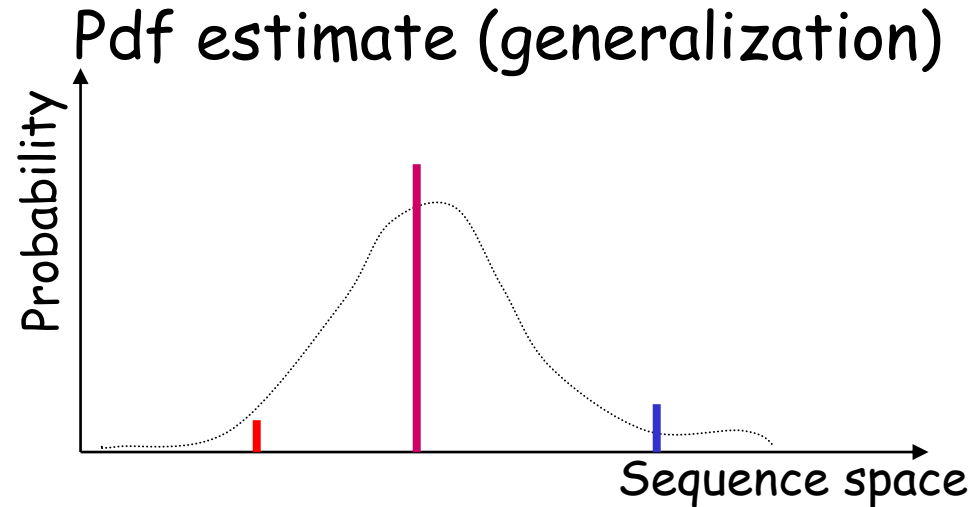
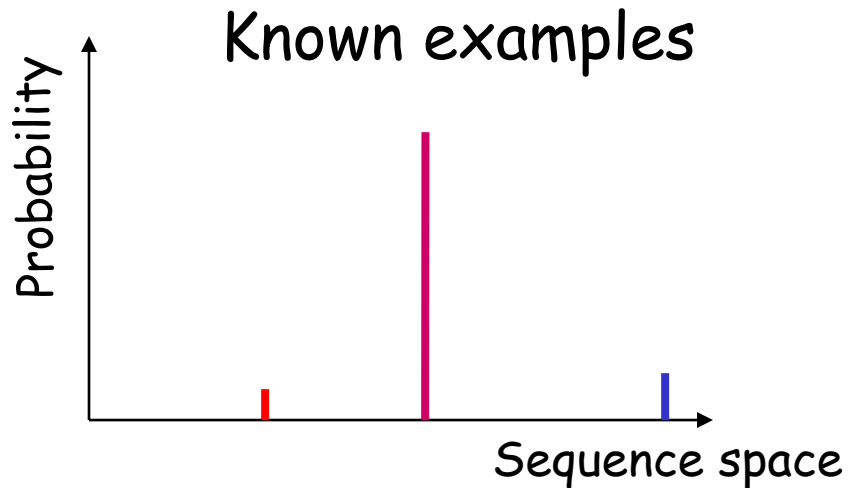
*Associates probability  $P(s_i | M)$  to  $s_i$*   
e.g.:  $M$  is the representation of the family of globins



# Probabilistic models for sequences

*Most useful probabilistic models are Trainable systems*

The probability density function over the sequence space is estimated from known examples by means of a learning algorithm



e.g.: Writing a generic representation of the sequences of globins starting from a set of known globins

# Probabilistic Models for Biological Sequences

- What are they?
- Why to use them?

# Modelling a protein family

## Probabilistic model

Given a protein class (e.g. Globins), a probabilistic model trained on this family can be adopted to compute a probability value for new sequences

Seq1	0.98
Seq2	0.21
Seq3	0.12
Seq4	0.89
Seq5	0.47
Seq6	0.78

This value measures the similarity between the new sequence and the family described by the model

# Probabilistic Models for Biological Sequences

- What are they?
- Why to use them?
- Which probabilities do they compute?

$P( s | M )$  or  $P( M | s )$  ?

A model  $M$  associates to a sequence  $s$  the probability  $P( s | M )$

This probability answers the question:

Which is the probability for a model  $M$  (e.g. describing the Globins) to generate the sequence  $s$  ?

The question we want to answer is:

Given a sequence  $s$ , does it belong to the class described by the model  $M$ ? (e.g. is it a Globin?)

We need to compute  $P( M | s )$  !!

# Bayes Theorem

$$P(X,Y) = P(X | Y) P(Y) = P(Y | X) P(X) \quad \text{Joint probability}$$

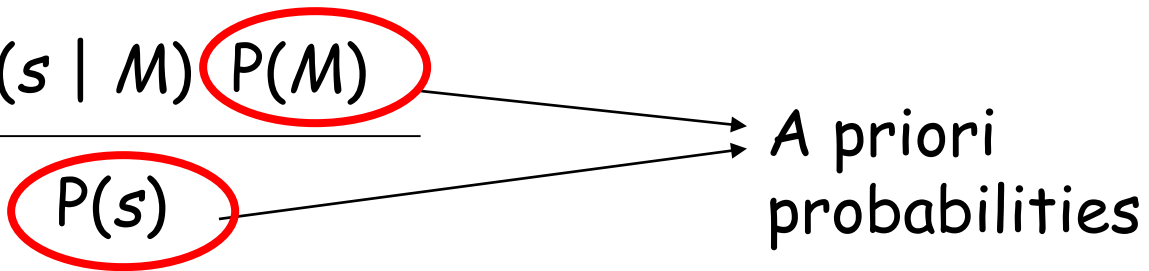


$$P(Y | X) = \frac{P(X | Y) P(Y)}{P(X)}$$

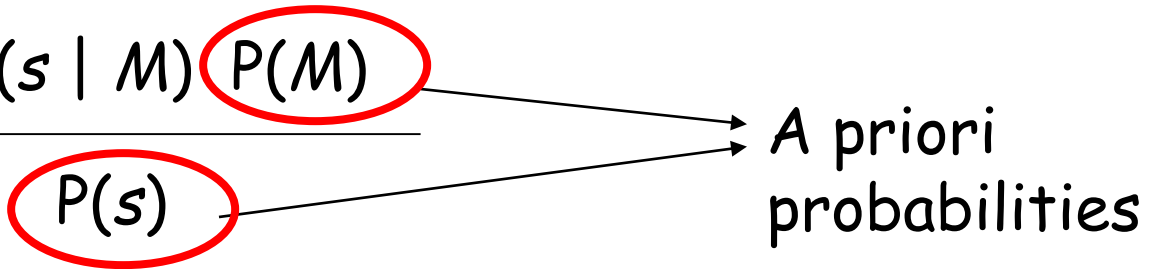
So:

$$P(M | s) = \frac{P(s | M) P(M)}{P(s)}$$

A priori probabilities

The terms P(M) in the numerator and P(s) in the denominator of the Bayes' theorem formula are each circled in red. Two black arrows originate from these circles and point towards the text 'A priori probabilities' on the right side of the slide.

## The *A priori* probabilities

$$P(M | s) = \frac{P(s | M) P(M)}{P(s)}$$


A priori probabilities

$P(M)$  is the probability of the model (i.e. of the class described by the model) BEFORE we know the sequence:

can be estimated as the abundance of the class

$P(s)$  is the probability of the sequence in the sequence space.

Cannot be reliably estimated!!

## Comparison between models

We can overcome the problem comparing the probability of generating  $s$  from different models

$$\frac{P(M_1 | s)}{P(M_2 | s)} = \frac{P(s | M_1) P(M_1)}{P(s)} \frac{P(s)}{P(s | M_2) P(M_2)} =$$

$$= \frac{P(s | M_1) P(M_1)}{P(s | M_2) P(M_2)}$$

Ratio between the abundances of the classes

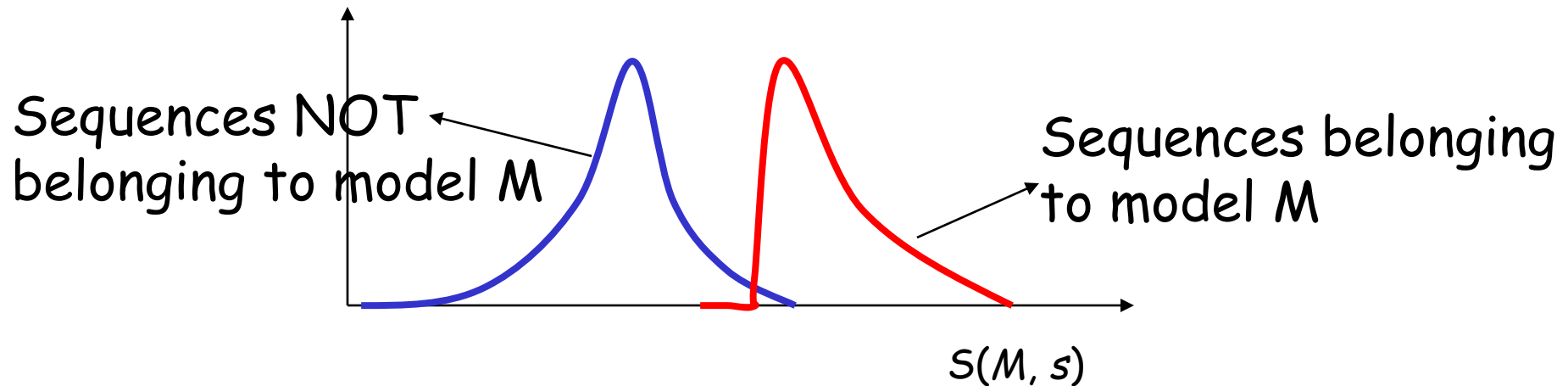


## Null model

Alternatively, we can score a sequence for a model  $M$  comparing it to a Null Model:

a model that generates ALL the possible sequences with probabilities depending ONLY on letter (e.g. residue) statistical abundance

$$S(M, s) = \log \frac{P(s | M)}{P(s | N)}$$



In this case we need a threshold and a statistic for evaluating the significance (E-value, P-value)

# The simplest probabilistic models: Markov Models

- Definition

## Markov Models Example: Weather

Register the weather conditions day by day:

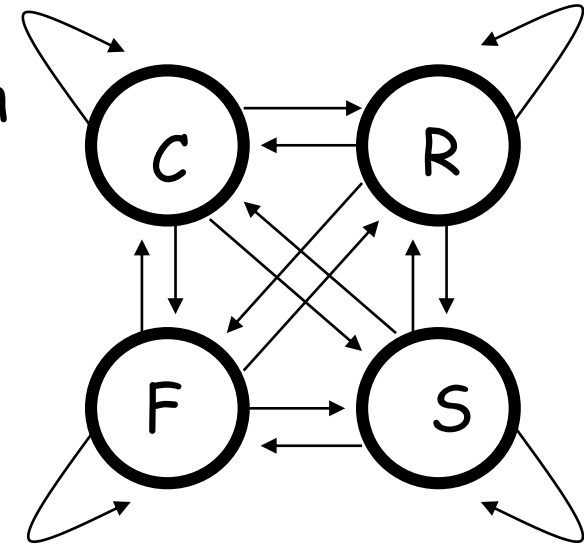
as a first hypothesis the weather condition in a day probabilistically depends ONLY on the weather conditions in the day before .

Define the conditional probabilities

$P(C|C)$ ,  $P(C|R)$ , ....  $P(R|C)$ .....

The probability for the 5-days registration  
CRRCS

$$P(CRRCS) = P(C) \cdot P(R|C) \cdot P(R|R) \cdot P(C|R) \cdot P(S|C)$$



C: Clouds

R: Rain

F: Fog

S: Sun

# Markov Model

Stochastic generator of sequences in which the probability of state in position  $i$  depends ONLY on the state in position  $i-1$

Given a set of states (== alphabet)  $\mathcal{C} = \{c_1; c_2; c_3; \dots c_N\}$

a Markov model is described with  $N \times (N+2)$  parameters

$\{a_{rt}, a_{\text{BEGIN } t}, a_{r \text{ END}}; r, t \in \mathcal{C}\}$

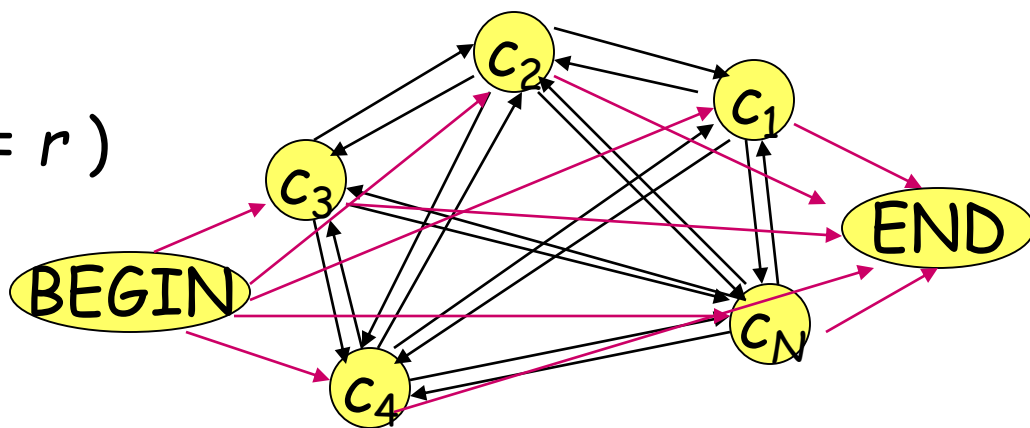
$$a_{rq} = P(s^i = q \mid s^{i-1} = r)$$

$$a_{\text{BEGIN } q} = P(s^1 = q)$$

$$a_{r \text{ END}} = P(s^T = \text{END} \mid s^{T-1} = r)$$

$$\sum_t a_{rt} + a_{r \text{ END}} = 1 \quad \forall r$$

$$\sum_t a_{\text{BEGIN } t} = 1$$



# Markov Models

Given the sequence:

$$s = s^1 s^2 s^3 s^4 s^6 \dots s^T$$

with  $s^i \in \mathcal{C} = \{c_1; c_2; c_3; \dots c_N\}$

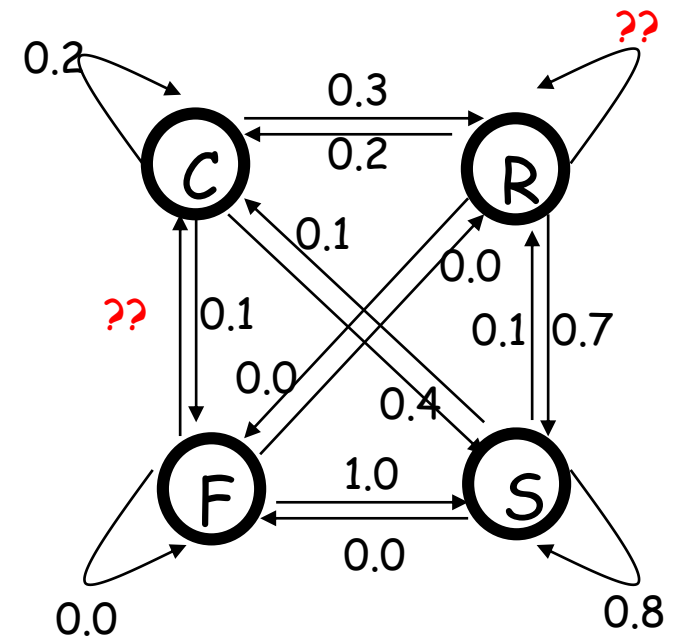
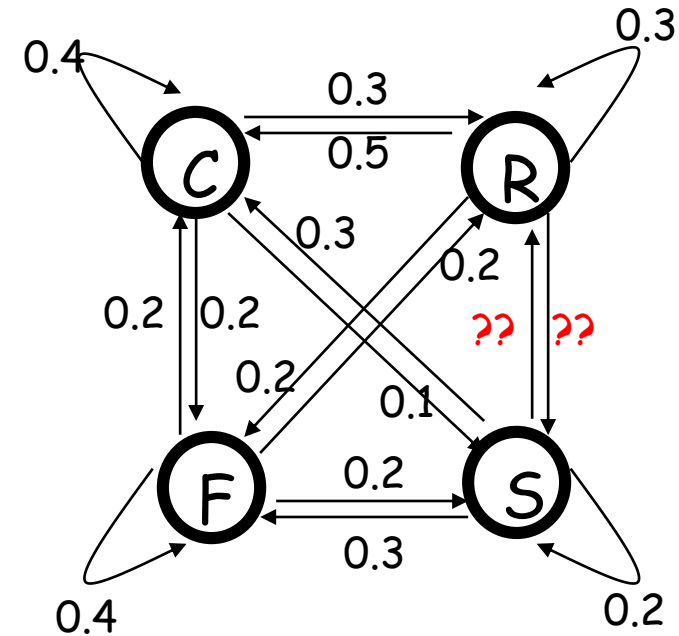
$$P(s \mid M) = P(s^1) \cdot \prod_{i=2} P(s^i \mid s^{i-1}) =$$

$$= a_{\text{BEGIN } s^1} \cdot \prod_{i=2}^T a_{s^{i-1} s^i} \cdot a_{s^T \text{ END}}$$

$$P(\text{"ALKALI"}) = a_{\text{BEGIN A}} \cdot a_{A L} \cdot a_{L K} \cdot a_{K A} \cdot a_{A L} \cdot a_{L I} \cdot a_{I \text{ END}}$$

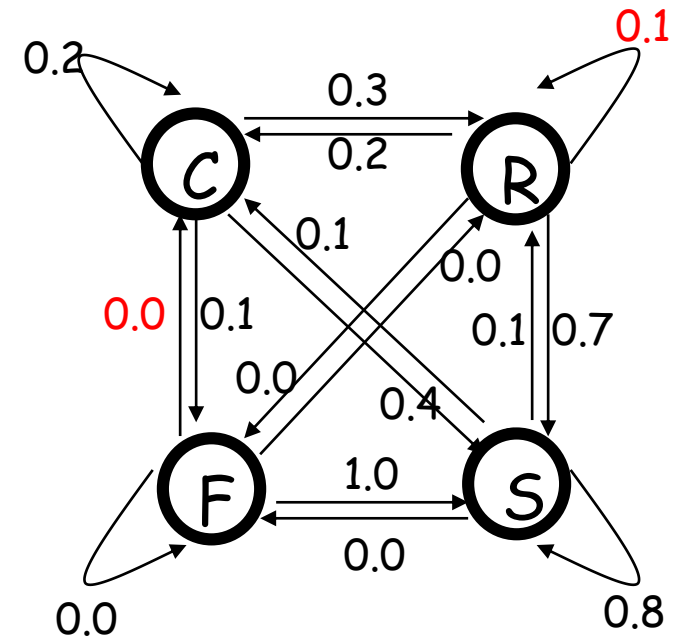
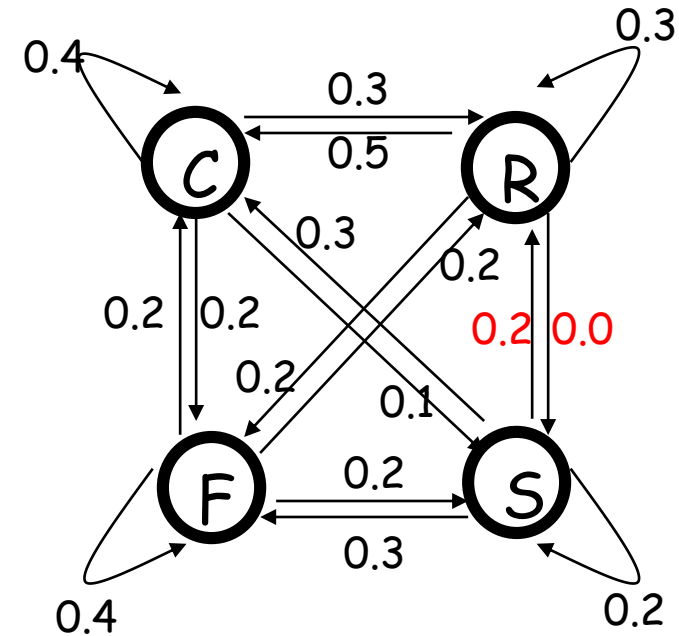
## Markov Models: Exercise

1) Fill the non defined values for the transition probabilities



## Markov Models: Exercise

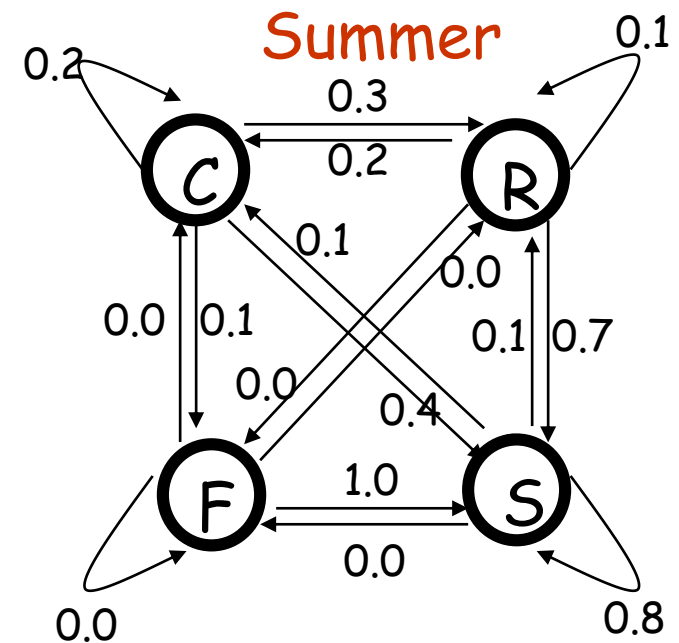
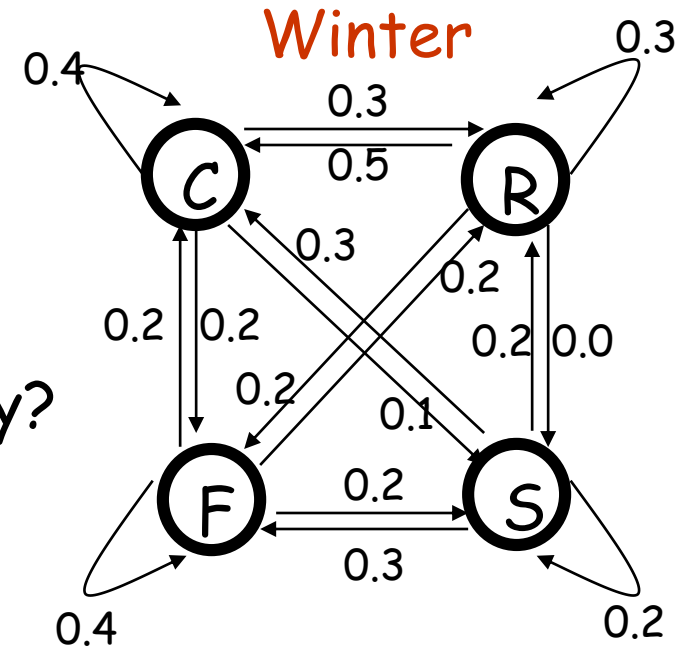
2) Which model better describes the weather in summer? Which one better describes the weather in winter?



## Markov Models: Exercise

3) Given the sequence  
CSSSCFS

which model gives the higher probability?  
[Consider the starting probabilities:  
 $P(X|\text{BEGIN})=0.25$ ]



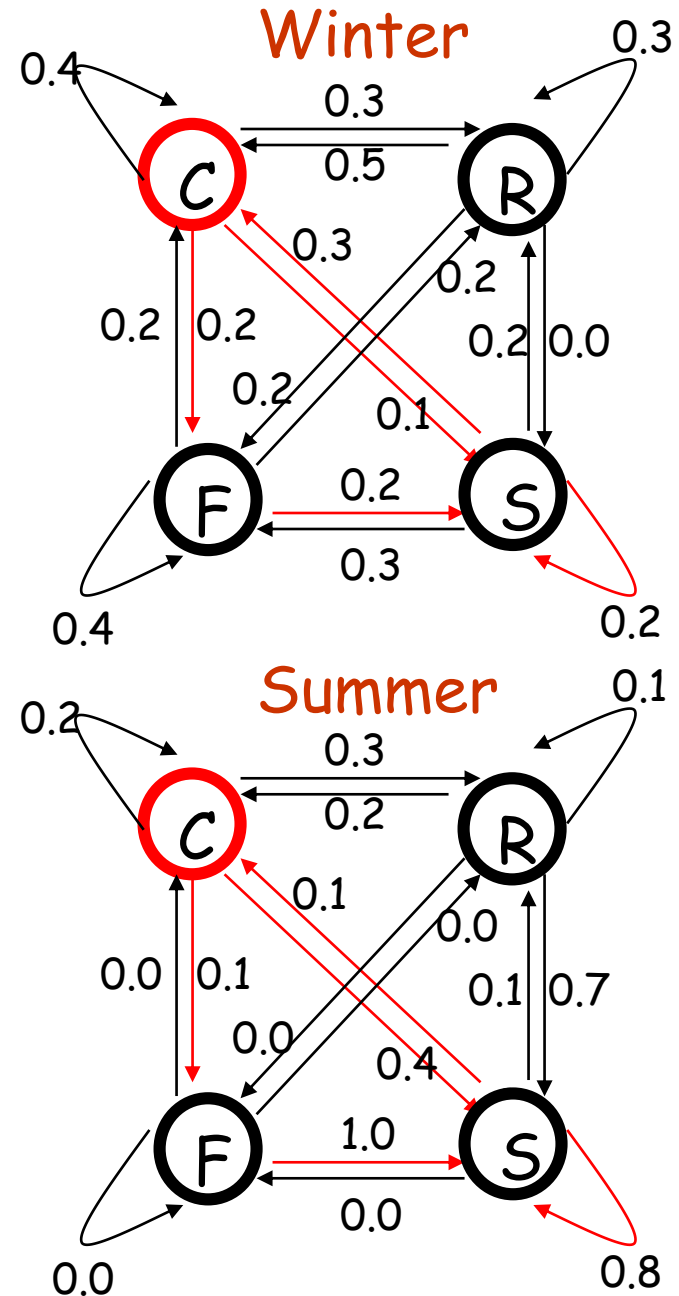


## Markov Models: Exercise

$$\begin{aligned} P(CSSSCFS \mid \text{Winter}) &= \\ &= 0.25 \times 0.1 \times 0.2 \times 0.2 \times 0.3 \times 0.2 \times 0.2 = \\ &= 1.2 \times 10^{-5} \end{aligned}$$

$$\begin{aligned} P(CSSSCFS \mid \text{Summer}) &= \\ &= 0.25 \times 0.4 \times 0.8 \times 0.8 \times 0.1 \times 0.1 \times 1.0 = \\ &= 6.4 \times 10^{-4} \end{aligned}$$

4) Can we conclude that the observation sequence refers to a summer week?

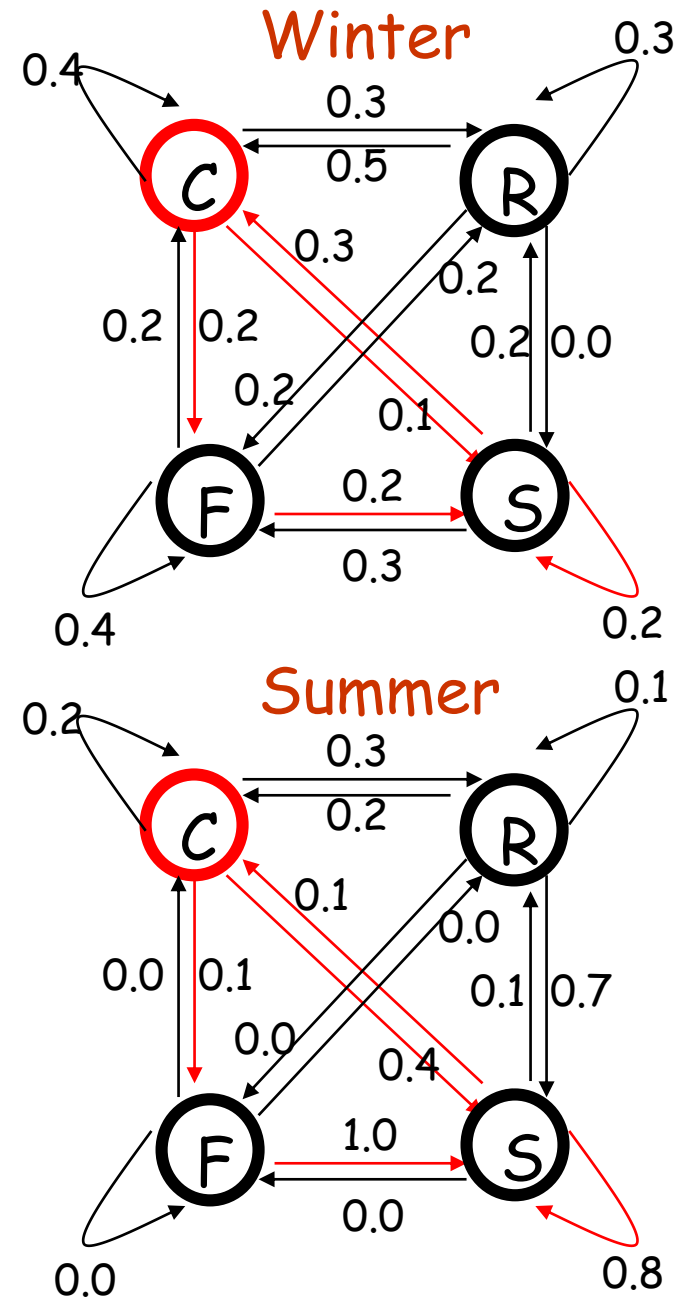


## Markov Models: Exercise

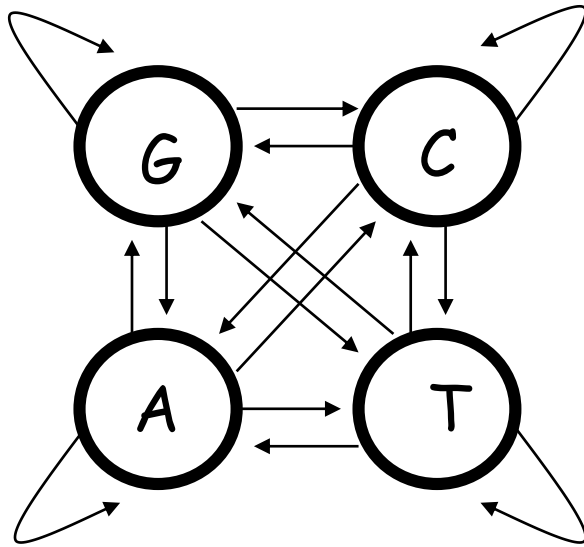
$$P(\text{Seq} \mid \text{Winter}) = 1.2 \times 10^{-5}$$

$$P(\text{Seq} \mid \text{Summer}) = 6.4 \times 10^{-4}$$

$$\frac{P(\text{Summer} \mid \text{Seq})}{P(\text{Winter} \mid \text{Seq})} = \frac{P(\text{Seq} \mid \text{Summer}) P(\text{Summer})}{P(\text{Seq} \mid \text{Winter}) P(\text{Winter})}$$



# Simple Markov Model for DNA sequences



DNA:

$C = \{\text{Adenine, Cytosine, Guanine, Thymine}\}$

16 transition probabilities (12 of which independent) +

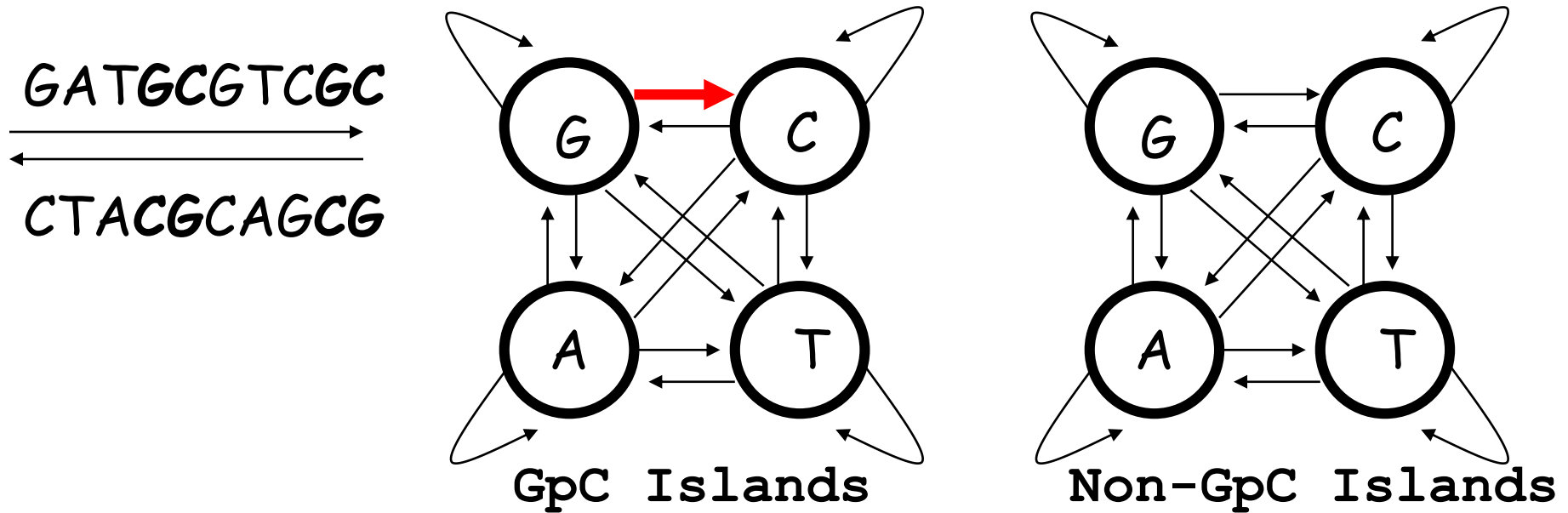
4 Begin probabilities +

4 End probabilities.

The parameters of the model are different in different zones of DNA

They describe the overall composition and the couple recurrences

# Example of Markov Models: GpC Island



In the Markov Model of GpC Islands  $a_{GC}$  is higher than in Markov Model Non-GpC Islands

Given a sequence  $s$  we can evaluate

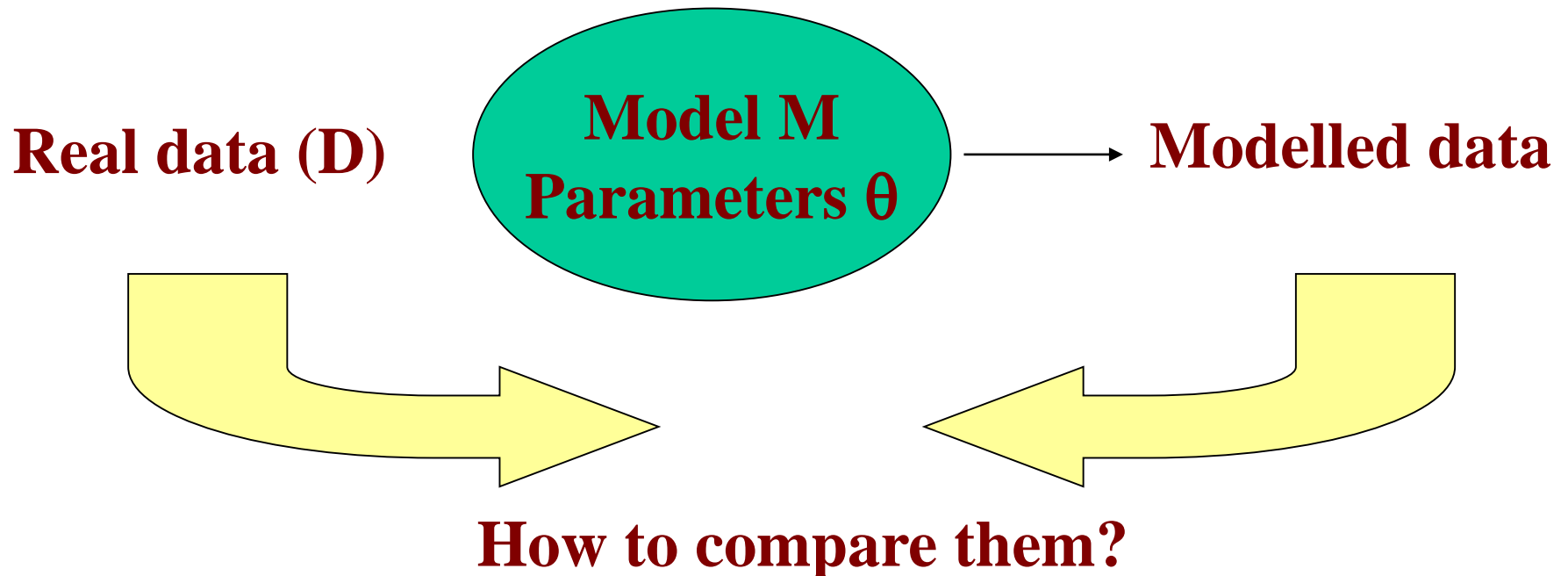
$$P(\text{GpC} \mid s) = \frac{P(s \mid \text{GpC}) \cdot P(\text{GpC})}{P(s \mid \text{GpC}) \cdot P(\text{GpC}) + P(s \mid \text{nonGpC}) \cdot P(\text{nonGpC})}$$

# The simplest probabilistic models: Markov Models

- Definition
- Training

# Probabilistic training of a parametric method

Generally speaking, a parametric model  $M$  aims to reproduce a set of known data



# Training of Markov Models

Let  $\theta_M$  be the set of parameters of model  $M$ .

During the training phase,  $\theta_M$  parameters are estimated from the set of known data  $\mathbf{D}$

*Maximum Likelihood Estimation (ML)*

$$\theta^{ML} = \operatorname{argmax}_{\theta} P(\mathbf{D} \mid M, \theta)$$

## Maximum Likelihood training: Proof

Given a sequence  $s$  contained in  $\mathbf{D}$ :

$$s = s^1 s^2 s^3 s^4 s^6 \dots\dots\dots s^T$$

$$P(s | M) = a_{BEGIN, s^1} \cdot \prod_{i=2}^{T-1} a_{s^i s^{i+1}} \cdot a_{s^T END}$$

We can count the number of transitions between any to states  $j$  and  $k$ :  $n_{jk}$

$$P(s | M) = \prod_{j=0}^{N+1} \prod_{k=0}^{N+1} a_{jk}^{n_{jk}} \quad \text{Where states 0 and } N+1 \text{ are BEGIN and END}$$

On top of this, keep in mind that normalisation constraints must be satisfied for each state

$$\forall j: \sum_{k'=0}^N a_{jk'} = 1$$

So the likelihood has to be maximised on the variety defined by the normalisation constraints. How to do that?



## Maximum Likelihood training: Proof

Given a sequence  $s$  contained in  $\mathbf{D}$ :

$$s = s^1 s^2 s^3 s^4 s^6 \dots\dots\dots s^T$$

$$P(s | M) = a_{BEGIN, s^1} \cdot \prod_{i=2}^{T-1} a_{s^i s^{i+1}} \cdot a_{s^T END}$$

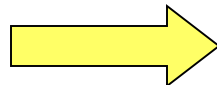
We can count the number of transitions between any to states  $j$  and  $k$ :  $n_{jk}$

$$P(s | M) = \prod_{j=0}^{N+1} \prod_{k=0}^{N+1} a_{jk}^{n_{jk}} \quad \text{Where states 0 and } N+1 \text{ are BEGIN and END}$$

Normalisation constraints are taken into account using the Lagrange multipliers  $\lambda_k$

$$L(\alpha_{ab}, \lambda_1, \lambda_2, \dots, \lambda_N) = P(s | M) - \sum_{j=0}^N \lambda_j \cdot \left( \sum_{k=0}^N a_{jk} - 1 \right)$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial a_{jk}} = \frac{n_{jk}}{a_{jk}} P(s | M) - \lambda_j = 0 \\ \frac{\partial L}{\partial \lambda_j} = \sum_{k'=0}^N a_{jk'} - 1 = 0 \end{array} \right.$$



$$a_{jk} = \frac{n_{jk}}{\sum_{k'=0}^N n_{jk'}}$$

# Training of Markov Models

Let  $\theta_M$  be the set of parameters of model  $M$ .

During the training phase,  $\theta_M$  parameters are estimated from the set of known data  $\mathbf{D}$

## *Maximum Likelihood Estimation (ML)*

$$\theta^{ML} = \operatorname{argmax}_{\theta} P( \mathbf{D} / M, \theta )$$

It can be proved that:

$$a_{ik} = \frac{n_{ik}}{\sum_j n_{ij}} \quad \begin{array}{l} \text{Frequency of occurrence as counted in the} \\ \text{data set } \mathbf{D} \end{array}$$

## *Maximum A Posteriori Estimation (MAP)*

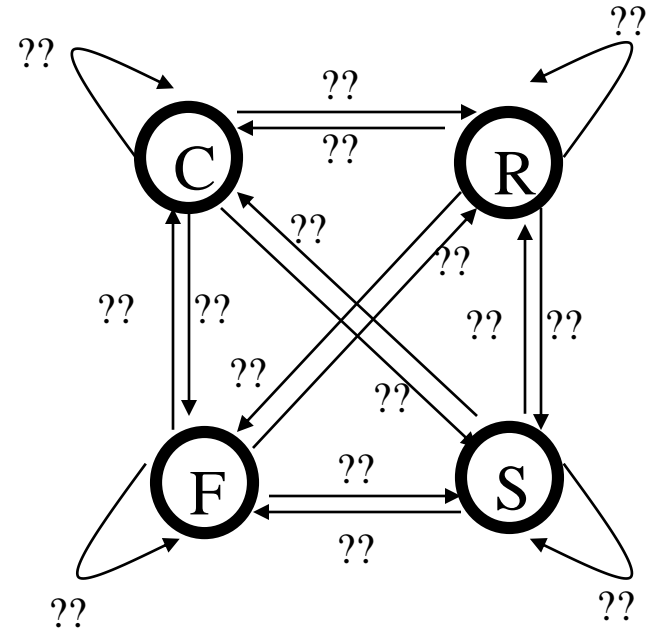
$$\theta^{MAP} = \operatorname{argmax}_{\theta} P( \theta / M, \mathbf{D} ) = \operatorname{argmax}_{\theta} [ P( \mathbf{D} / M, \theta ) \cdot P(\theta) ]$$

# Training of Markov Models: Exercise

Given the observation sequence:

CCCFFCRRRRCCSSSSFSFRRFFSSF

set the parameters of the Markov Model.



# Hidden Markov Models

- Preliminary examples

## Loaded dice

We have 99 regular dice ( $R$ ) and 1 loaded die ( $L$ ).

	$P(1)$	$P(2)$	$P(3)$	$P(4)$	$P(5)$	$P(6)$
$R$	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$
$L$	$1/10$	$1/10$	$1/10$	$1/10$	$1/10$	$1/2$

Given a sequence:

4156266656321636543662152611536264162364261664616263

We don't know the sequence of dice that generated it.

RRRRRLRLRRRRRRRLRRRRRRRRRRRRRLRLRRRRRRRRRLRRRRLRRRRLRR

# Loaded dice

*Hypothesis:*

We chose a different die for each roll

Two stochastic processes give origin to the sequence of observations.

- 1) Choosing the die (  $R$  o  $L$  ).
- 2) Rolling the die

The sequence of dice is *hidden*

The first process is assumed to be Markovian (in this case a 0-order MM)

The outcome of the second process depends only on the state reached in the first process (that is the chosen die)

# Casinò

## Model



Each state ( $R$  and  $L$ ) generates a character of the alphabet  $\mathcal{C} = \{1, 2, 3, 4, 5, 6\}$

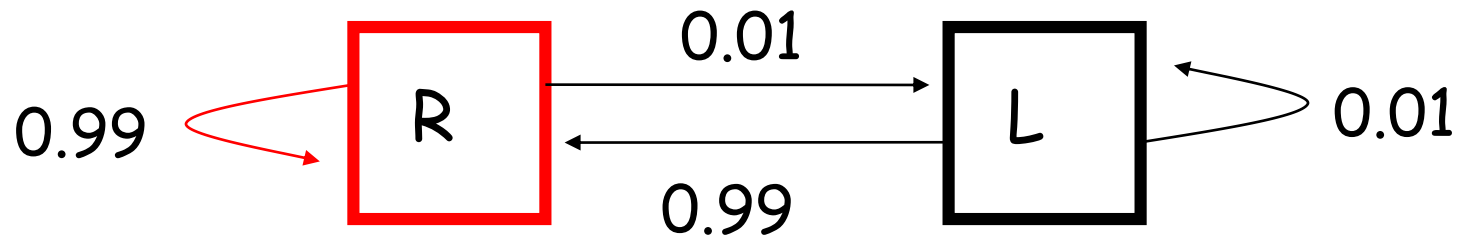
The *emission probabilities* depend only on the state.

The *transition probabilities* describe a Markov model that generates a state path: the *hidden sequence* ( $\pi$ )

The *observations sequence* ( $s$ ) is generated by two concomitant stochastic processes

*The observations sequence (s) is generated by two concomitant stochastic processes*

415626665632163654366215261  
RRRRRLRLRRRRRRRLRRRRRRRRRRRR



**Choose the State : R**

**Probability= 0.99**

**Chose the Symbol: 1**

**Probability= 1/6 (given R)**

415626665632163654366215261**1**  
RRRRRLRLRRRRRRRLRRRRRRRRRRRR**R**



*The observations sequence (s) is generated by two concomitant stochastic processes*

4156266656321636543662152611  
RRRRRLRLRRRRRRRLRRRRRRRRRRRR



**Choose the State : L**

**Probability= 0.99**

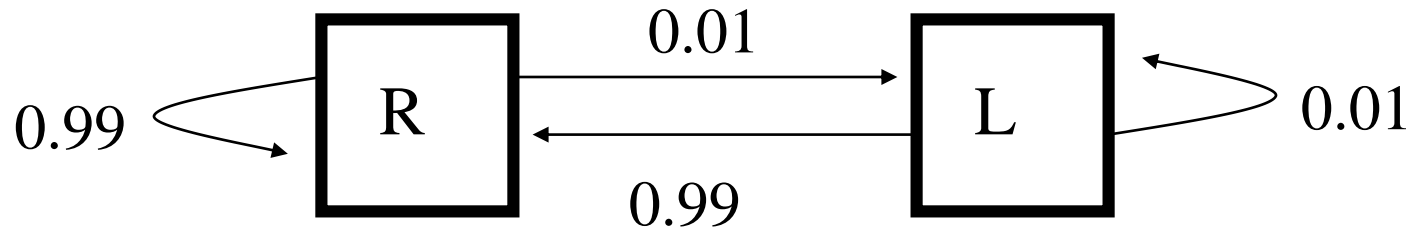
**Chose the Symbol: 5**

**Probability= 1/10 (given L)**

4156266656321636543662152611**5**  
RRRRRLRLRRRRRRRLRRRRRRRRRRRR**L**

## Loaded dice

*Model*



Each state (**R** and **L**) generates a character of the alphabet  
 $C = \{1, 2, 3, 4, 5, 6\}$

The *emission probabilities* depend only on the state.

The *transition probabilities* describe a Markov model that generates a state path: the hidden sequence ( $\pi$ )

The *observations sequence* ( $s$ ) is generated by two concomitant stochastic processes

# Why to search for the hidden path?

## Some not so serious example

### 1) DEMOGRAPHY

*Observable:* Number of births and deaths in a year in a village.

*Hidden variable:* Economic conditions (as a first approximation we can consider the success in business as a random variable, and by consequence, the wealth as a Markov variable

---> can we deduce the economic conditions of a village during a century by means of the register of births and deaths?

### 2) THE METEOROPATHIC TEACHER

*Observable:* Average of the marks that a meteoropathic teacher gives to their students during a day.

*Hidden variable:* Weather conditions

---> can we deduce the weather conditions during a years by means of the class register?

# Why to search for the hidden path?

To be more serious

## 1) SECONDARY STRUCTURE

*Observable:* protein sequence

*Hidden variable:* secondary structure

---> can we deduce (predict) the secondary structure of a protein given its amino acid sequence?

## 2) ALIGNMENT

*Observable:* protein sequence

*Hidden variable:* position of each residue along the alignment of a protein family

---> can we align a protein to a family, starting from its amino acid sequence?

# Hidden Markov Models

- Preliminary examples
- Formal definition

# Formal definition of Hidden Markov Models

A HMM is a stochastic generator of sequences characterised by:

- $N$  states
- A set of transition probabilities between two states  $\{a_{kj}\}$   
 $a_{kj} = P(\pi(i) = j / \pi(i-1) = k)$
- A set of starting probabilities  $\{a_{0k}\}$   
 $a_{0k} = P(\pi(1) = k)$
- A set of ending probabilities  $\{a_{k0}\}$   
 $a_{k0} = P(\pi(i) = \text{END} / \pi(i-1) = k)$
- An alphabet  $C$  with  $M$  characters.
- A set of emission probabilities for each state  $\{e_k(c)\}$

$$e_k(c) = P(s^i = c / \pi(i) = k)$$

- Constraints:

$$\sum_k a_{0k} = 1$$

$$a_{k0} + \sum_j a_{kj} = 1 \quad \forall k$$

$$\sum_{c \in C} e_k(c) = 1 \quad \forall k$$

$s$ : sequence

$\pi$ : path through the states

# Generating a sequence with a HMM

Choose the initial state  $\pi(1)$  following the probabilities  $a_{0k}$

$i = 1$

Choose the character  $s^i$  from the alphabet  $\mathcal{C}$  following the probabilities  $e_k(c)$

Choose the next state following the probabilities  $a_{kj}$  and  $a_{k0}$

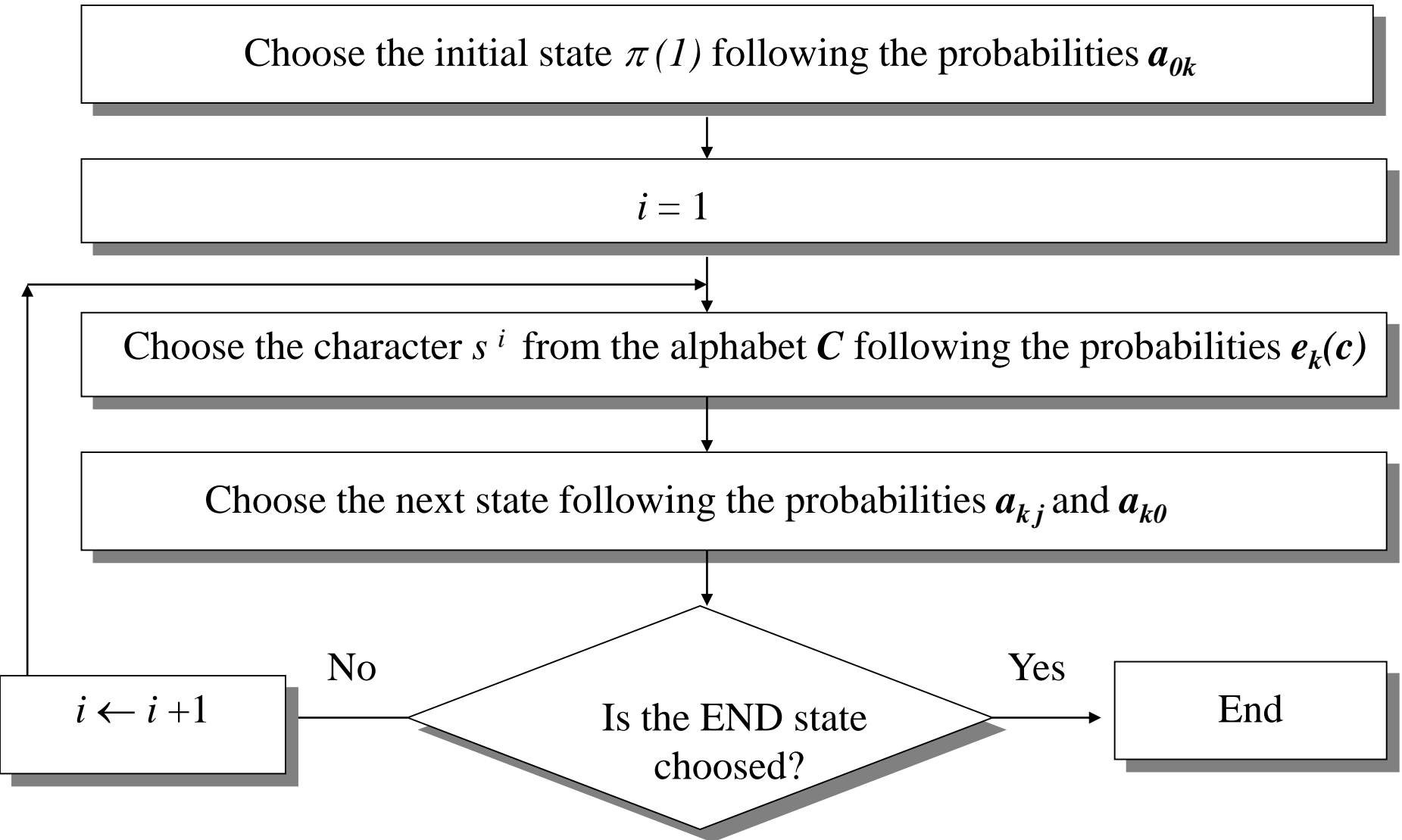
Is the END state  
chosen?

No

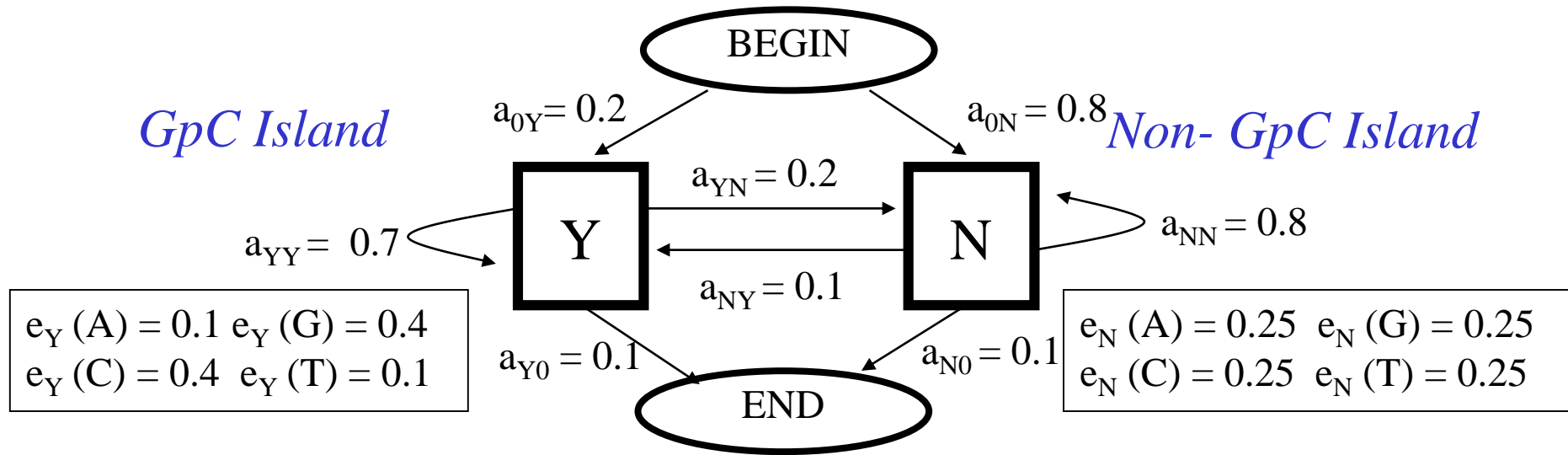
$i \leftarrow i + 1$

Yes

End



## GpC Island, simple model



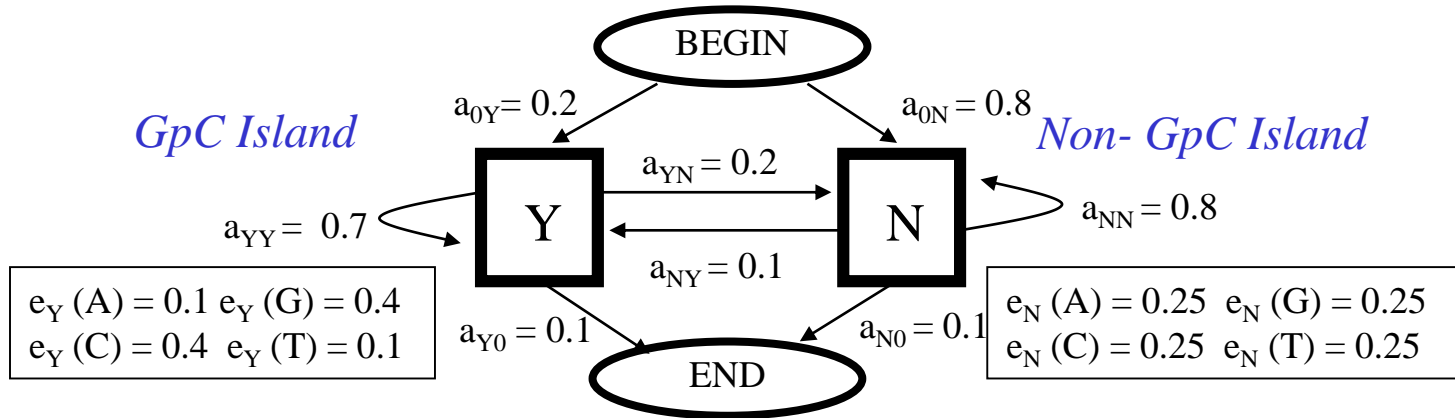
$s$  : **AGCGCGTAATCTG**

$\pi$  : **YYYYYYNNNNNN**

•  $P(s, \pi / M)$  can be easily computed



## $P(s, \pi | M)$ can be easily computed



$s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \quad \mathbf{G} \quad \mathbf{C} \quad \mathbf{G} \quad \mathbf{T} \quad \mathbf{A} \quad \mathbf{A} \quad \mathbf{T} \quad \mathbf{C} \quad \mathbf{T} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{Y} \quad \mathbf{N} \quad \mathbf{N} \quad \mathbf{N} \quad \mathbf{N} \quad \mathbf{N} \quad \mathbf{N}$

Emission:  $0.1 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.1 \times 0.25 \times 0.25 \times 0.25 \times 0.25 \times 0.25 \times 0.25$   
 Transition:  $0.2 \times 0.7 \times 0.7 \times 0.7 \times 0.7 \times 0.7 \times 0.7 \times 0.2 \times 0.8 \times 0.8 \times 0.8 \times 0.8 \times 0.8 \times 0.1$

Multiplying all the probabilities gives the probability of having the sequence AND the path through the states

## Evaluation of the joint probability of the sequence and the path

$$P(s, \pi \mid M) = P(s \mid \pi, M) \cdot P(\pi \mid M)$$

$$P(\pi \mid M) = a_{0\pi(1)} \cdot \prod_{i=2}^T a_{\pi(i-1)\pi(i)} \cdot a_{\pi(T)0}$$

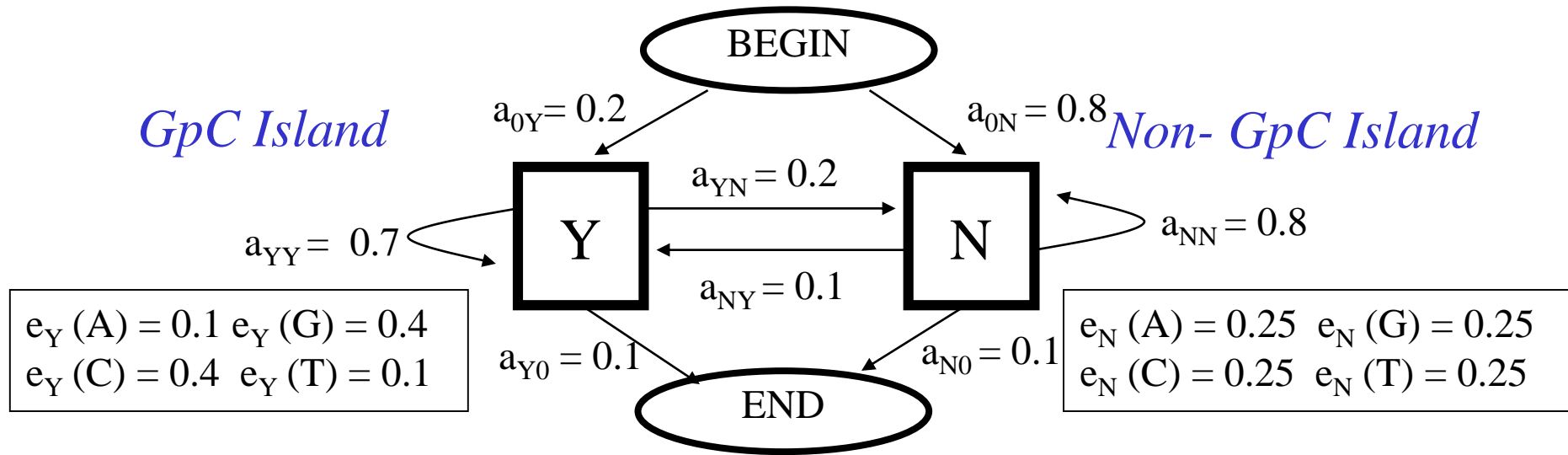
$$P(s \mid \pi, M) = \prod_{i=1}^T e_{\pi(i)}(s^i)$$

$$P(s, \pi \mid M) = a_{\pi(T)0} \cdot \prod_{i=1}^T a_{\pi(i-1)\pi(i)} \cdot e_{\pi(i)}(s^i)$$

# Hidden Markov Models

- Preliminary examples
- Formal definition
- Three questions

## GpC Island, simple model



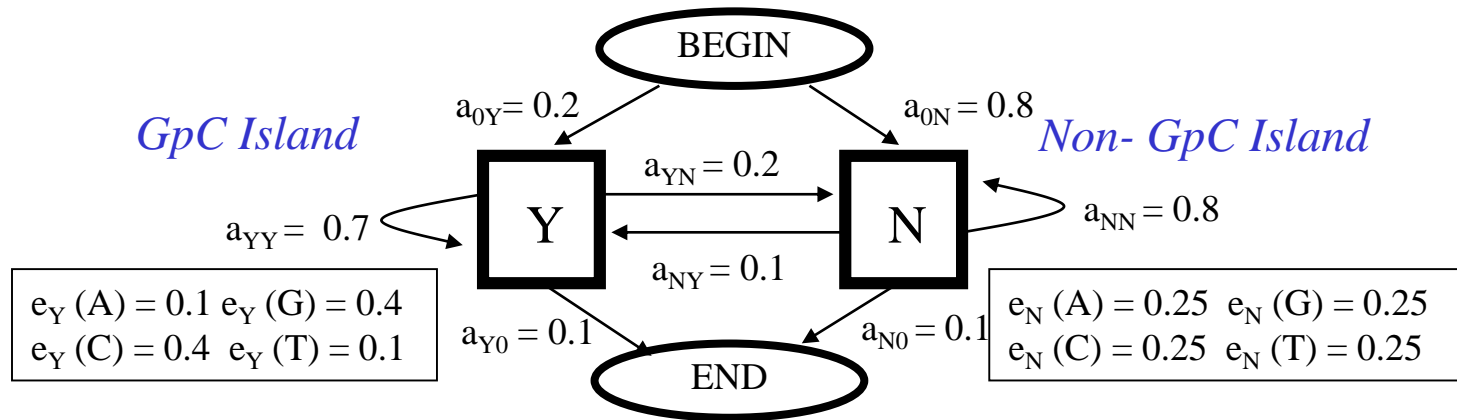
$s$  : **AGCGCGTAATCTG**

$\pi$ : ?????????????

•  $P(s, \pi / M)$  can be easily computed

• *How to evaluate  $P(s / M)$ , when the path is unknown?*

# How to evaluate $P(s | M)$ ?



$$P(s | M) = \sum_{\pi} P(s, \pi | M)$$

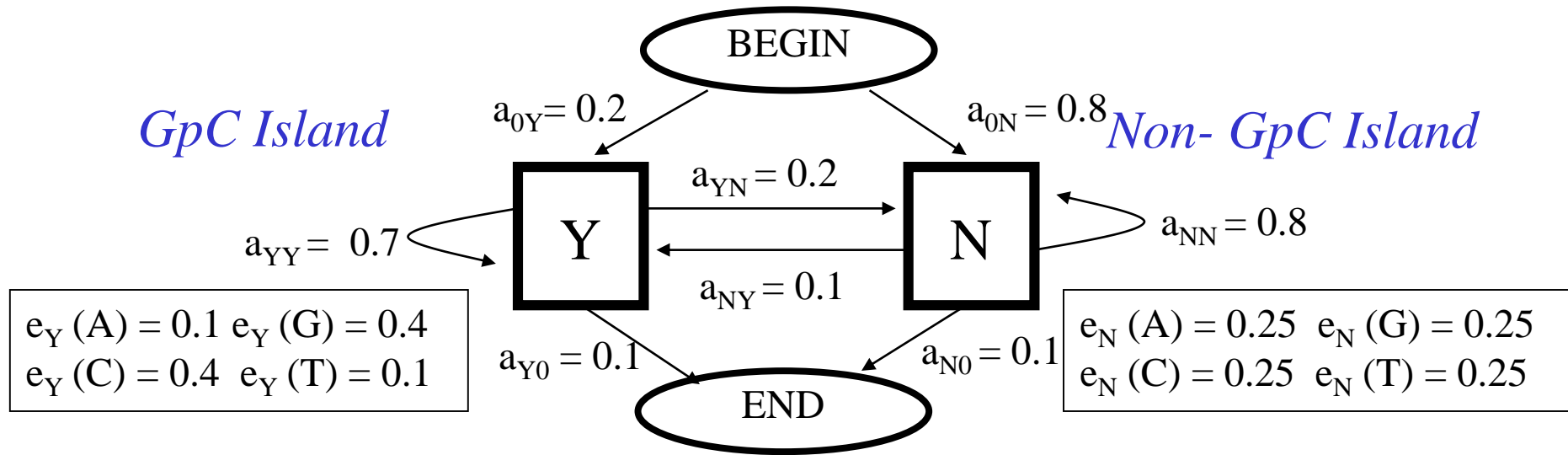
$s :$	<b>A</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>T</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>G</b>
$\pi_1 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>
$\pi_2 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>
$\pi_3 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>
$\pi_4 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>N</b>
$\pi_5 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>	<b>Y</b>

.....

**$2^{13}$  different paths**

Summing over all the path will give the probability of having the sequence

## Resumé: GpC Island, simple model

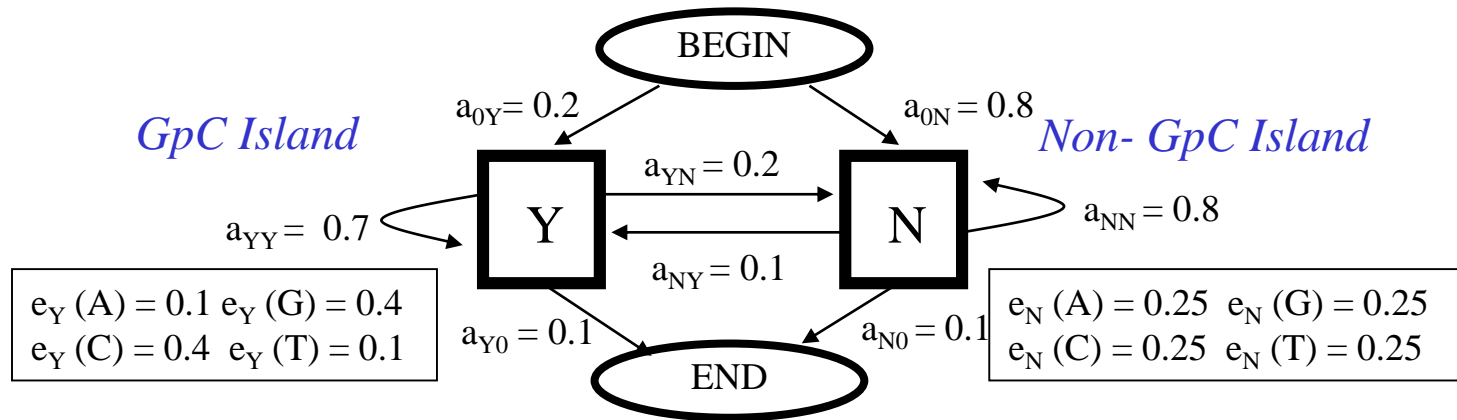


$s$  : **AGCGCGTAATCTG**

$\pi$  : **????????????**

- $P(s, \pi / M)$  can be easily computed
- How to evaluate  $P(s / M)$ , when the path is unknown?
- Can we show the hidden path?

## Can we show the hidden path?



$$\pi^* = \operatorname{argmax}_{\pi} [ P( \pi | s, M ) ] = \operatorname{argmax}_{\pi} [ P( \pi, s | M ) ]$$

$s$ :	<b>A</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>T</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>G</b>
$\pi_1$ :	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>
$\pi_2$ :	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>
$\pi_3$ :	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>
$\pi_4$ :	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>N</b>
$\pi_5$ :	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>	<b>Y</b>

.....

**$2^{13}$  different paths**

**Viterbi path: path that gives the best joint probability**

## Can we show the hidden path?

### *Viterbi decoding*

Among all the possible path, choose the path  $\pi^*$  that maximises the  $P(\pi | s, M)$

$$\pi^* = \operatorname{argmax}_{\pi} [ P(\pi | s, M) ] = \operatorname{argmax}_{\pi} [ P(\pi, s | M) ]$$

### *A Posteriori decoding*

For each position choose the state  $\underline{\pi}(t)$  :

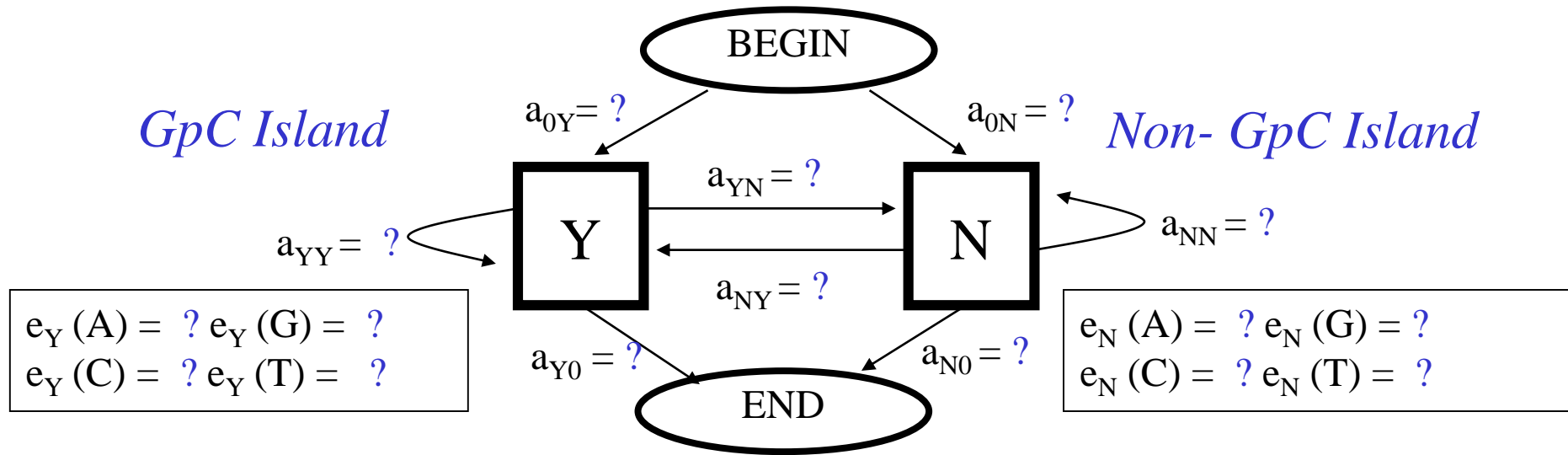
$$\underline{\pi}(i) = \operatorname{argmax}_k [ P(\pi(i) = k | s, M) ]$$

The contribution to this probability derives from all the paths that go through the state  $k$  at position  $i$ .

The A posteriori path can be a non-sense path (it may not be a legitimate path if some transitions are not permitted in the model)



## GpC Island, simple model

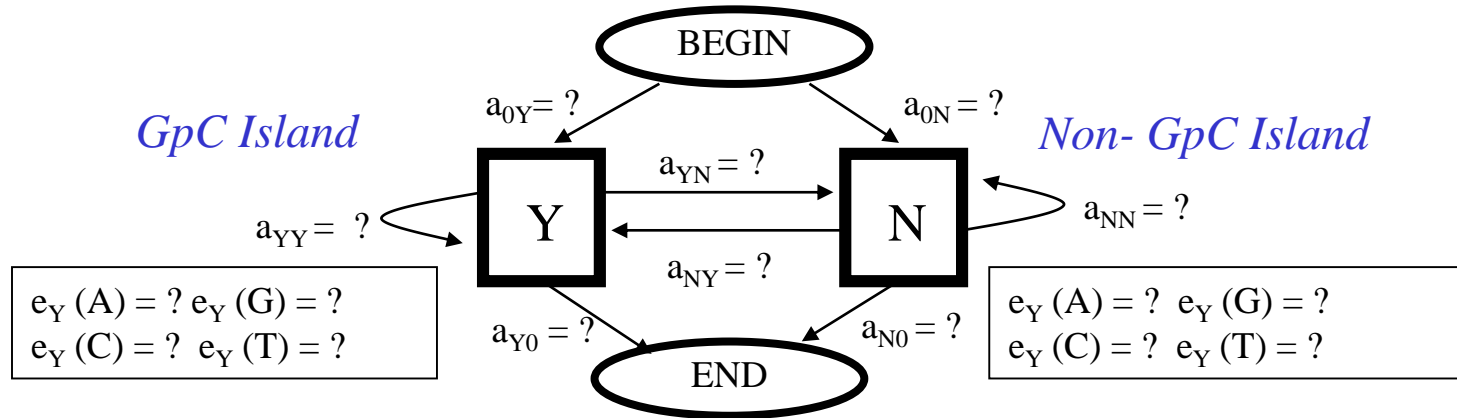


$s$  : **AGCGCGTAATCTG**

$\pi$  : **YYYYYYNNNNNN**

- $P(s, \pi / M)$  can be easily computed
- How to evaluate  $P(s / M)$ , when the path is unknown?
- Can we show the hidden path?
- Can we evaluate the parameters starting from known examples?

# Can we evaluate the parameters starting from known examples?



$S : \mathbf{A \quad G \quad C \quad G \quad C \quad G \quad T \quad A \quad A \quad T \quad C \quad T \quad G}$

$\pi : \mathbf{Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad N \quad N \quad N \quad N \quad N \quad N}$

Emission:  $e_Y(A) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(T) \times e_N(A) \times e_N(A) \times e_N(T) \times e_N(C) \times e_N(T) \times e_N(G)$

Transition:  $a_{0Y} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{N0}$

How to find the parameters  $e$  and  $a$  that maximises this probability?

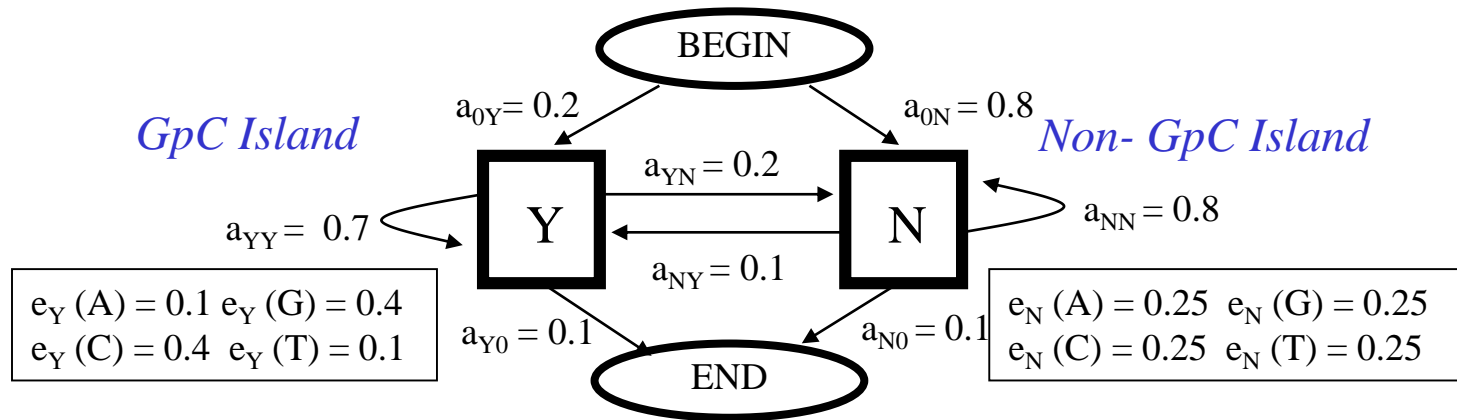
How if we don't know the path?

# Hidden Markov Models: Algorithms

- Résumé
- Evaluating  $P(s \mid M)$ : Forward Algorithm



# Strategy: Summing over all the possible paths



$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{Y} \quad \mathbf{Y}$

Emission:  $0.1 \times 0.4$

Transition:  $0.2 \times 0.7$  *0.0056*

$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{N} \quad \mathbf{Y}$

Emission:  $0.25 \times 0.4$

Transition:  $0.8 \times 0.1$  *0.008*

$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{Y} \quad \mathbf{N}$

Emission:  $0.1 \times 0.25$

Transition:  $0.2 \times 0.2$  *0.001*

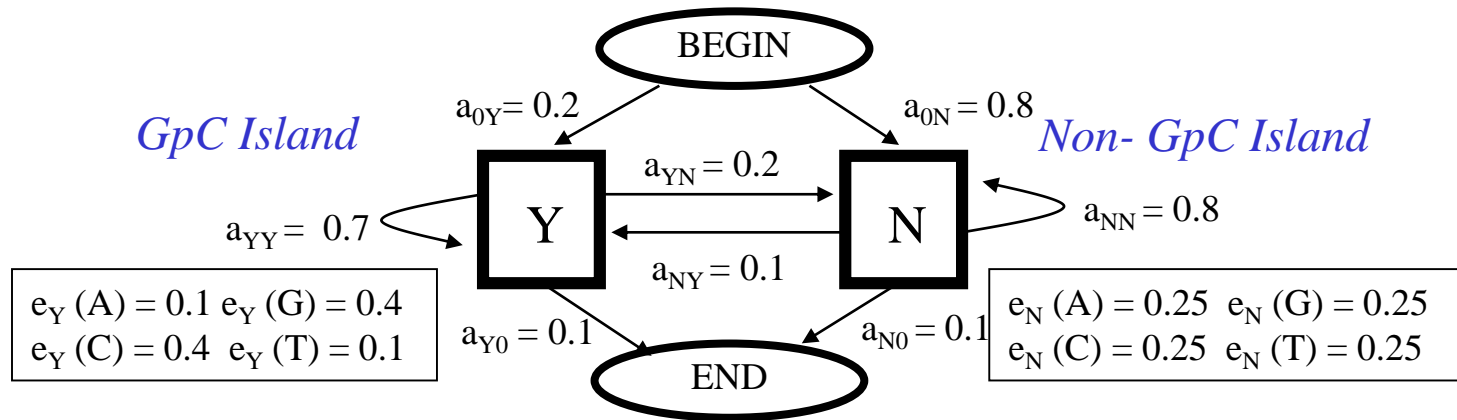
$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{N} \quad \mathbf{N}$

Emission:  $0.25 \times 0.25$

Transition:  $0.8 \times 0.8$  *0.04*

# Strategy: Summing over all the possible paths



$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{Y} \quad \mathbf{Y}$

Emission:  $0.1 \times 0.4$

Transition:  $0.2 \times 0.7$  *0.0056*

$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{N} \quad \mathbf{Y}$

Emission:  $0.25 \times 0.4$

Transition:  $0.8 \times 0.1$  *0.008*

**In order to continue the process we only need to know the last state:**

$s : \mathbf{A} \quad \mathbf{G}$

$\pi : \mathbf{Y} \quad \mathbf{N}$

Emission:  $0.1 \times 0.25$

Transition:  $0.2 \times 0.2$  *0.001*

$s : \mathbf{A} \quad \mathbf{G}$

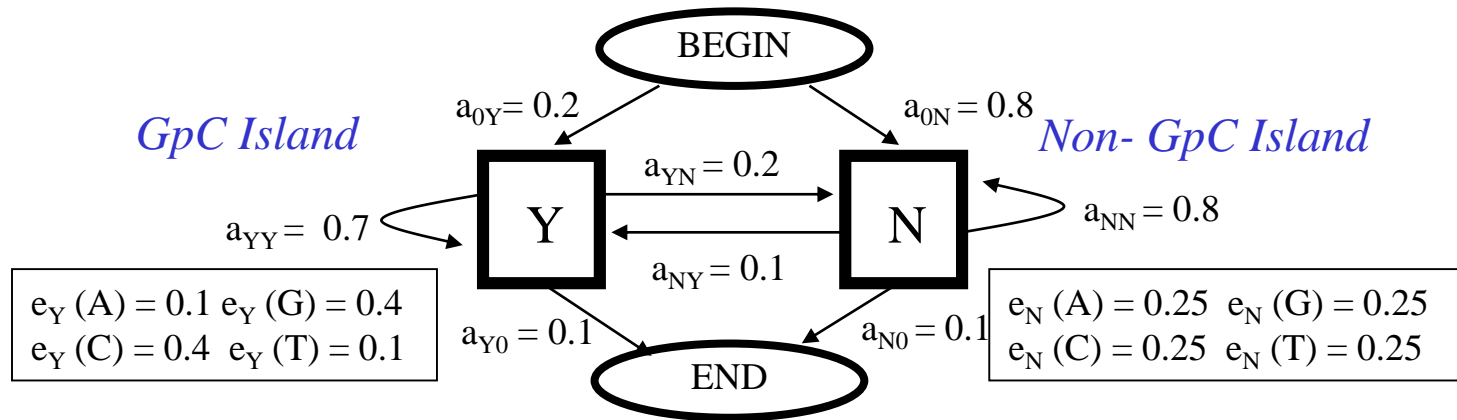
$\pi : \mathbf{N} \quad \mathbf{N}$

Emission:  $0.25 \times 0.25$

Transition:  $0.8 \times 0.8$  *0.04*

**The contributions from previous states can be summed up**

# Strategy: Summing over all the possible paths



$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{Y}$

Emission:  $0.1 \times 0.4$

Transition:  $0.2 \times 0.7$  *0.0056*

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{Y}$

Emission:  $0.25 \times 0.4$

Transition:  $0.8 \times 0.1$  *0.008*

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{X} \quad \mathbf{Y}$

*0.0136*

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{N}$

Emission:  $0.1 \times 0.25$

Transition:  $0.2 \times 0.2$  *0.001*

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{N}$

Emission:  $0.25 \times 0.25$

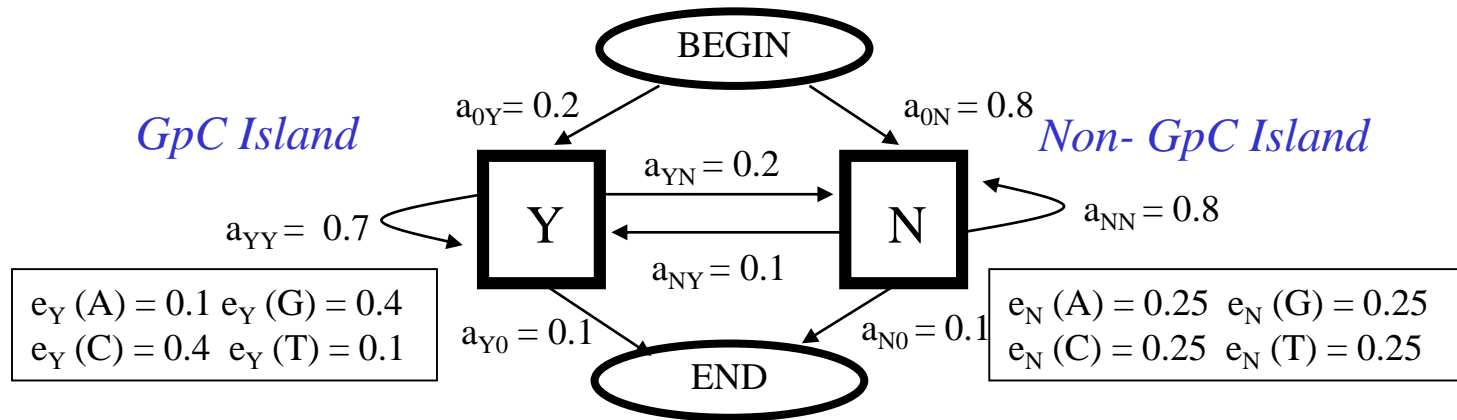
Transition:  $0.8 \times 0.8$  *0.04*

Sum

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{X} \quad \mathbf{N}$

*0.041*

# Strategy: Summing over all the possible paths

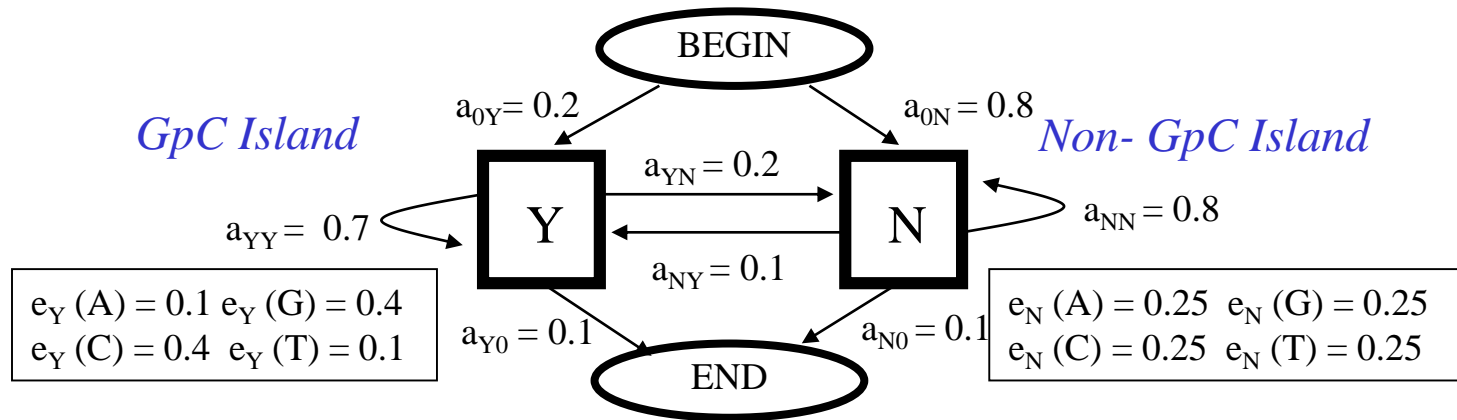


$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{Y} \quad \mathbf{Y} \\
 0.0136 \times \begin{array}{c} 0.4 \\ 0.7 \end{array} + 0.041 \times \begin{array}{c} 0.4 \\ 0.1 \end{array}
 \end{array}$$

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{Y} \quad \mathbf{N} \\
 0.0136 \times \begin{array}{c} 0.25 \\ 0.2 \end{array} + 0.041 \times \begin{array}{c} 0.25 \\ 0.8 \end{array}
 \end{array}$$



# Strategy: Summing over all the possible paths



$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{Y} \quad \mathbf{Y} \\
 0.0136 \times \begin{array}{c} 0.4 \\ 0.7 \end{array}
 \end{array}$$

+

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{N} \quad \mathbf{Y} \\
 0.041 \times \begin{array}{c} 0.4 \\ 0.1 \end{array}
 \end{array}$$

0.005448

Sum

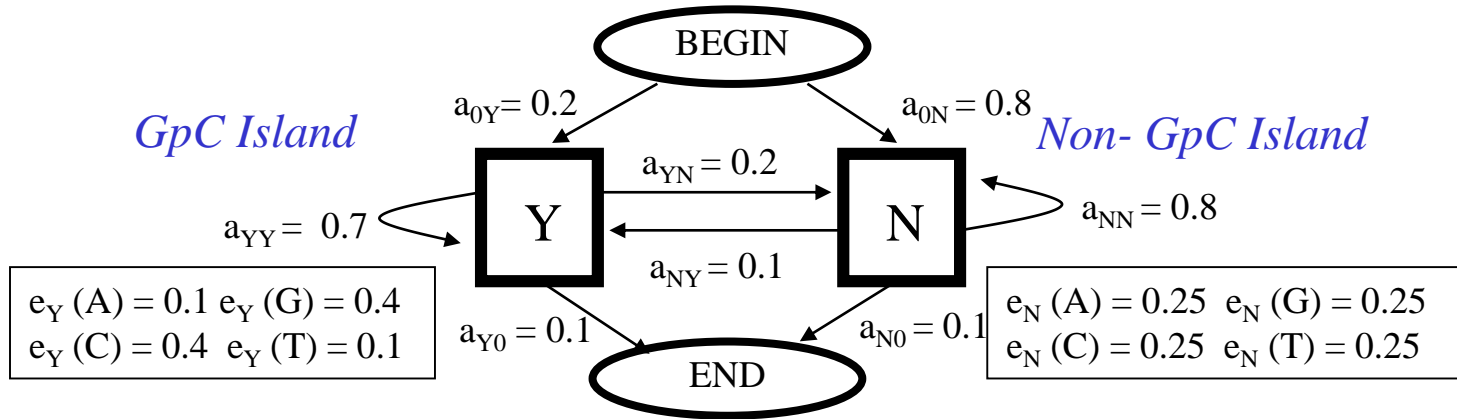
$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{Y} \quad \mathbf{N} \\
 0.0136 \times \begin{array}{c} 0.25 \\ 0.2 \end{array}
 \end{array}$$

+

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{X} \quad \mathbf{N} \quad \mathbf{N} \\
 0.041 \times \begin{array}{c} 0.25 \\ 0.8 \end{array}
 \end{array}$$

0.00888

## Strategy: Summing over all the possible paths



Iterating until the last position of the sequence:

<b>A</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>T</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>G</b>	<div> <div>×</div> <div>0.1 (a<sub>Y0</sub>)</div> <div>+</div> <div>×</div> <div>0.1 (a<sub>N0</sub>)</div> </div>
<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>Y</b>	
<b>A</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>T</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>G</b>	
<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>N</b>	
													P(s M)

## Forward Algorithm

On the basis of preceding observations the computation of  $P(s / M)$  can be decomposed in simplest problems

For each state  $k$  and each position  $i$  in the sequence, we compute:

$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k / M) \rightarrow \text{Will be understood}$$

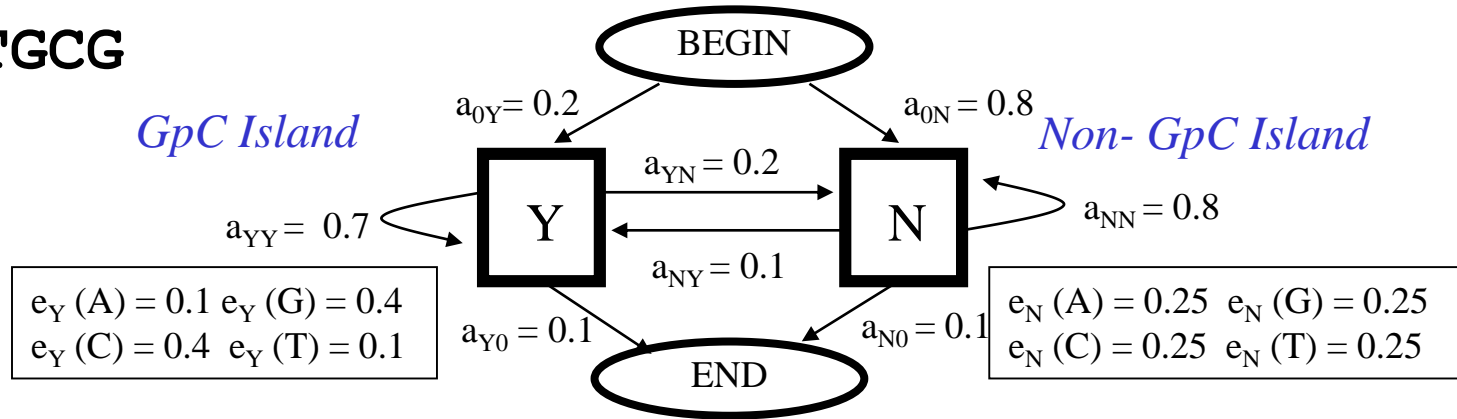
*Initialisation:*  $F_{BEGIN}(0) = 1 \quad F_i(0) = 0 \quad \forall i \neq BEGIN$

*Recurrence:* 
$$\begin{aligned} F_l(i+1) &= P(s^1 s^2 \dots s^i s^{i+1}, \pi(i+1) = l) = \\ &= \sum_k P(s^1 s^2 \dots s^i, \pi(i) = k) \cdot a_{kl} \cdot e_l(s^{i+1}) = \\ &= e_l(s^{i+1}) \cdot \sum_k F_k(i) \cdot a_{kl} \end{aligned}$$

*Termination:* 
$$\begin{aligned} P(s) &= P(s^1 s^2 s^3 \dots s^T, \pi(T+1) = END) = \\ &= \sum_k P(s^1 s^2 \dots s^T, \pi(T) = k) \cdot a_{k0} \\ &= \sum_k F_k(T) \cdot a_{k0} \end{aligned}$$

# Forward Algorithm: Example

**S = ATGCG**



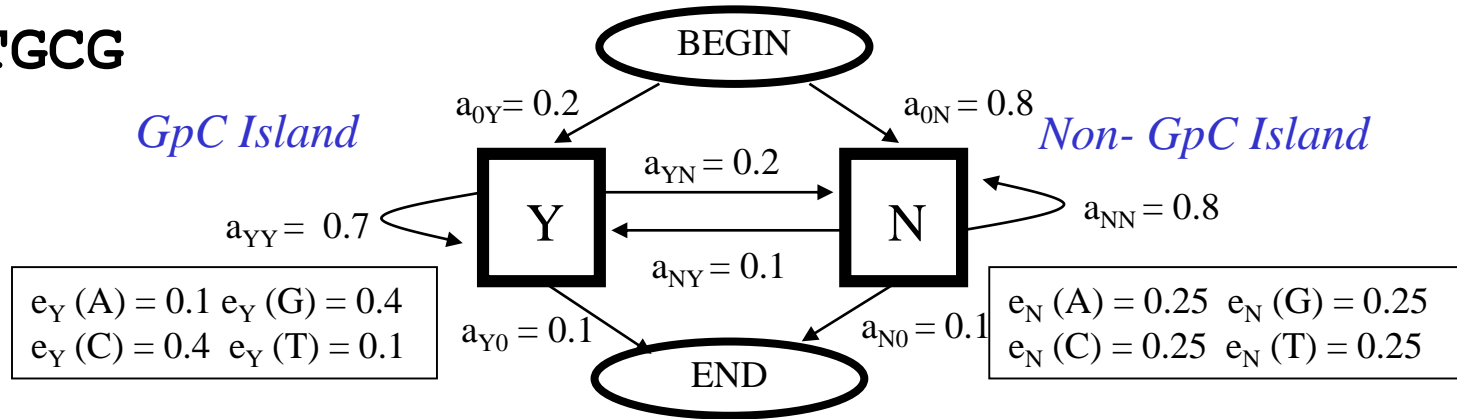
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Initialisation:*  $F_{BEGIN}(0) = 1$   $F_i(0) = 0 \quad \forall i \neq BEGIN$

	-	A	T	G	C	G	-
Begin	1						
Y	0						
N	0						
End	0						

# Forward Algorithm: Example

**S = ATGCG**



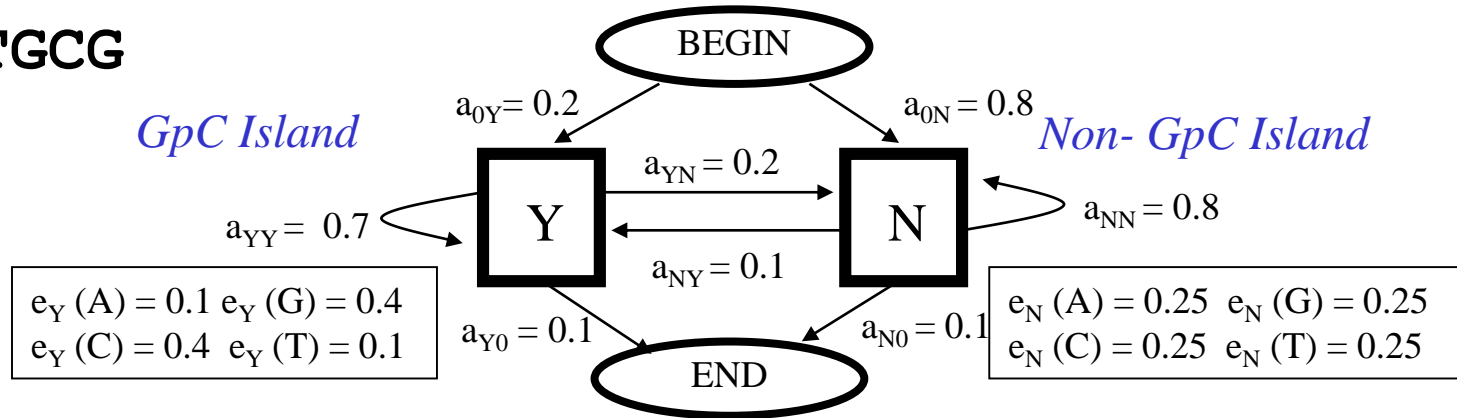
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Recurrence:*  $F_l(1) = e_l(s^1) \cdot \sum_k F_k(0) \cdot a_{kl}$

	-	A	T	G	C	G	-
Begin	1	0					
Y	0	$0.2 \times 0.1 = 0.02$					
N	0	$0.8 \times 0.25 = 0.2$					
End	0	0					

# Forward Algorithm: Example

**S = ATGCG**



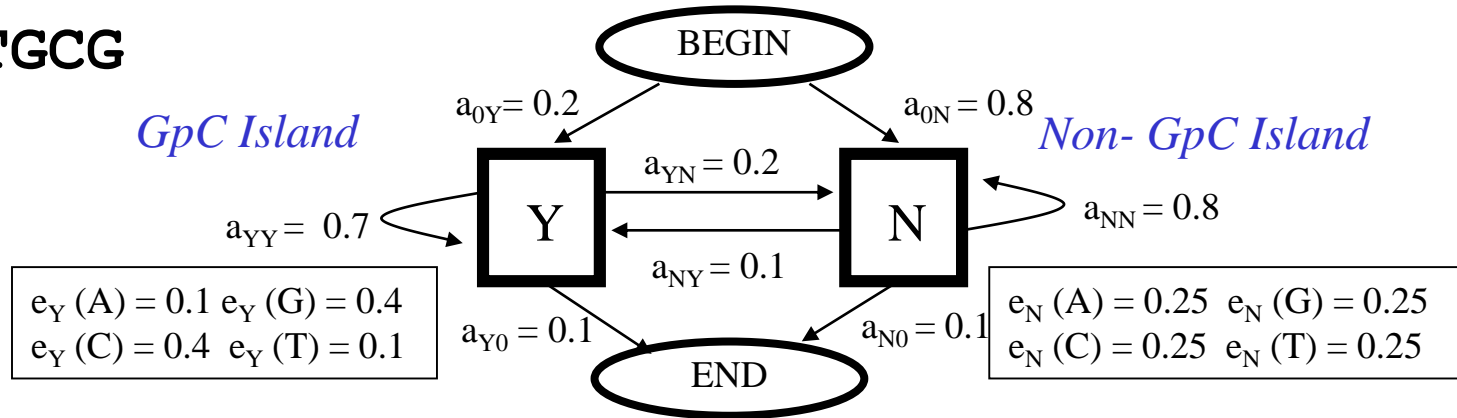
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Recurrence:*  $F_l(2) = e_l(s^2) \cdot \sum_k F_k(1) \cdot a_{kl}$

	-	A	T	G	C	G	-
Begin	1	0	0				
Y	0	2e-2	2e-2x0.7x0.1+ +0.2x0.1x0.1= <b>3.4e-3</b>				
N	0	0.2	2e-2x0.2x0.25+ +0.2x0.8x0.25= <b>4.1e-2</b>				
End	0	0	0				

# Forward Algorithm: Example

**S = ATGCG**



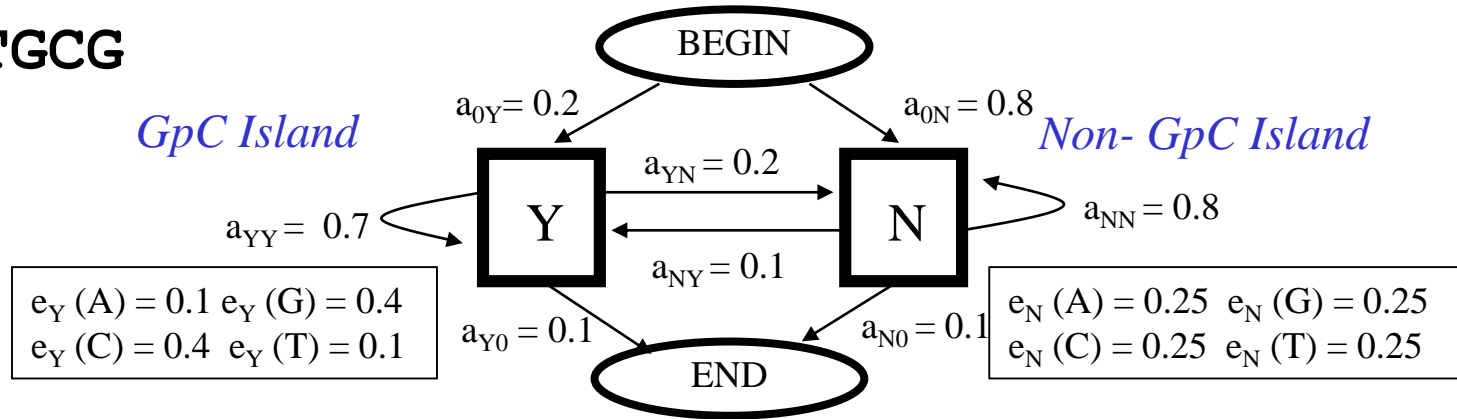
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Recurrence:*  $F_l(3) = e_l(s^3) \cdot \sum_k F_k(2) \cdot a_{kl}$

	-	A	T	G	C	G	-
Begin	1	0	0	0			
Y	0	2e-2	3.4e-3	3.4e-3x0.7x0.4+ +4.1e-2x0.1x0.4= <b>=2.592e-3</b>			
N	0	0.2	4.1e-2	3.4e-3x0.2x0.25+ +4.1e-2x0.8x0.25= <b>=8.37e-3</b>			
End	0	0	0	0			

# Forward Algorithm: Example

**S = ATGCG**



$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

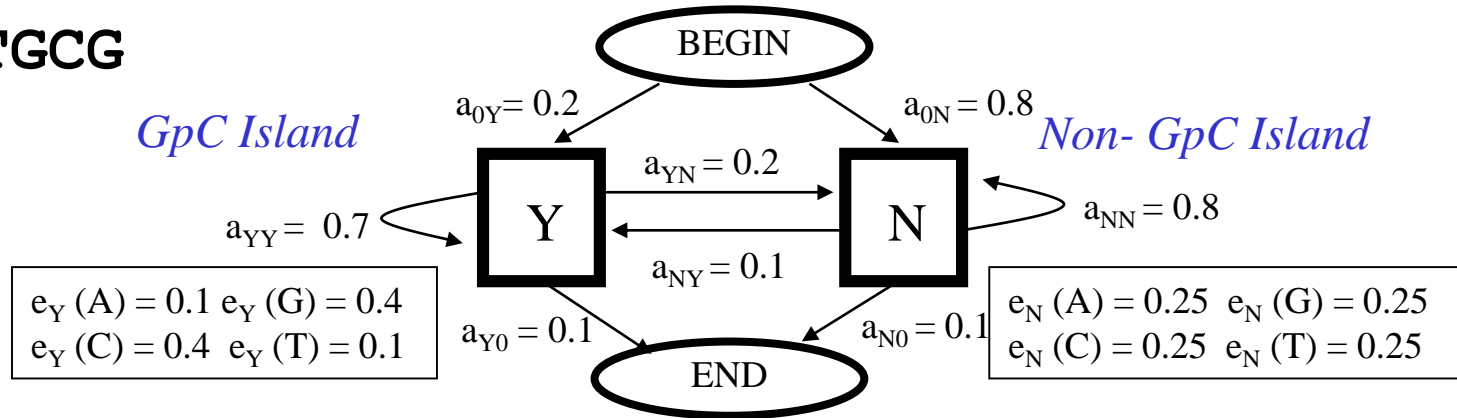
*Recurrence:*  $F_l(4) = e_l(s^4) \cdot \sum_k F_k(3) \cdot a_{kl}$

	-	A	T	G	C	G	-
Begin	1	0	0	0	0		
Y	0	2e-2	3.4e-3	2.592e-3	2.592e-3x0.7x0.4+ +8.37e-3x0.1x0.4= <b>=1.06056e-3</b>		
N	0	0.2	4.1e-2	8.37e-3	2.592e-3x0.2x0.25+ +8.37e-3x0.8x0.25= <b>=1.8036e-3</b>		
End	0	0	0	0	0		



# Forward Algorithm: Example

**S = ATGCG**



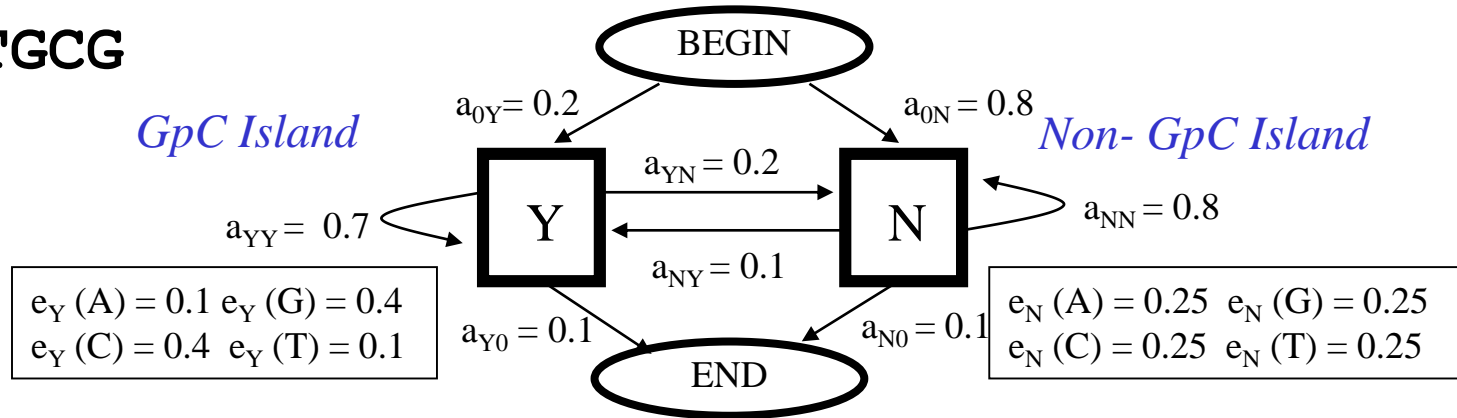
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Recurrence:*  $F_l(5) = e_l(s^5) \cdot \sum_k F_k(4) \cdot a_{kl}$

	-	A	T	G	C	G	-
Begin	1	0	0	0	0	0	
Y	0	2e-2	3.4e-3	2.592e-3	1.06056e-3	1.06056e-3x0.7x0.4+ +1.8036e-3x0.1x0.4= <b>=3.691008e-4</b>	
N	0	0.2	4.1e-2	8.37e-3	1.8036e-3	1.06056e-3x0.2x0.25+ +1.8036e-3x0.8x0.25= <b>=4.13748e-4</b>	
End	0	0	0	0	0	0	

# Forward Algorithm: Example

**S = ATGCG**



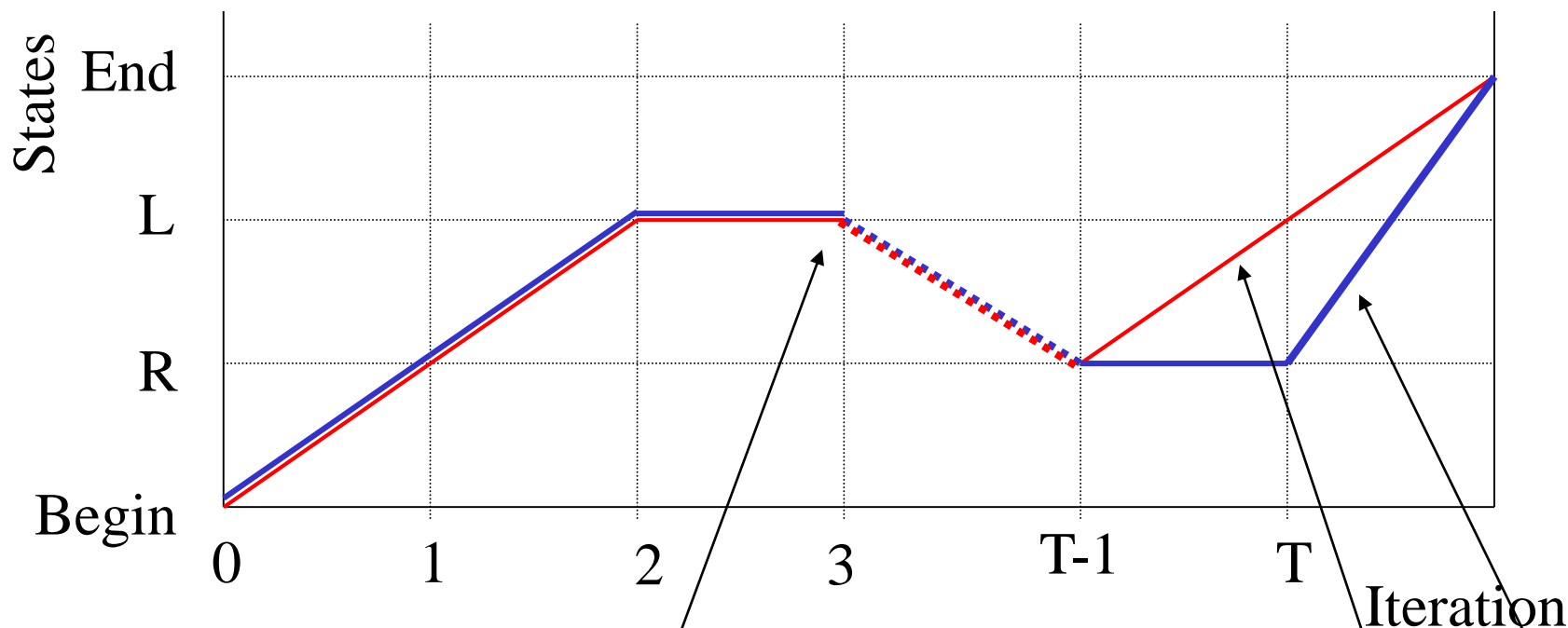
$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

*Termination:*  $P(s) = \sum_k F_k(T) \cdot a_{k0}$

	-	A	T	G	C	G	-
Begin	1	0	0	0	0	0	0
Y	0	2e-2	3.4e-3	2.592e-3	1.06056e-3	3.691008e-4	0
N	0	0.2	4.1e-2	8.37e-3	1.8036e-3	4.13748e-4	0
End	0	0	0	0	0	0	$3.691008e-4 \times 0.1 +$ $4.13748e-4 \times 0.1 =$ <b>7.828488e-5</b>

$$P(s|M) = 7.828488e-5$$

# Computing $P(s, \pi | M)$ for each path is a redundant operation

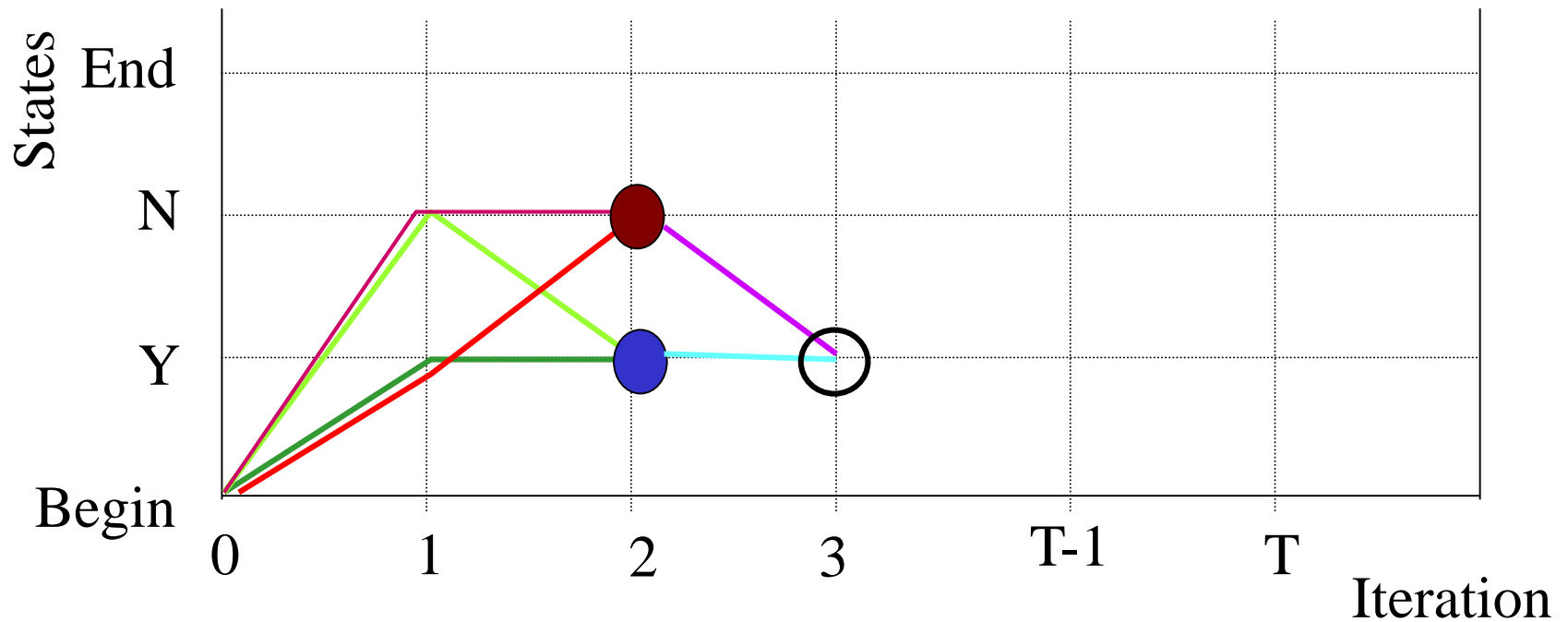


$$P(s, \pi_1 | M) = \left( \prod_{t=1}^{T-1} a_{\pi_1(t-1)\pi_1(t)} \cdot e_{\pi_1(t)}(s^t) \right) \cdot a_{\pi_1(T-1)\pi_1(T)} \cdot e_{\pi_1(T)}(s^T) \cdot a_{\pi_1(T)0}$$

$$P(s, \pi_2 | M) = \left( \prod_{t=1}^{T-1} a_{\pi_2(t-1)\pi_2(t)} \cdot e_{\pi_2(t)}(s^t) \right) \cdot a_{\pi_2(T-1)\pi_2(T)} \cdot e_{\pi_2(T)}(s^T) \cdot a_{\pi_2(T)0}$$

If we compute the common part only once we gain  $2 \cdot (T-1)$  operations

## Summing over all the possible paths



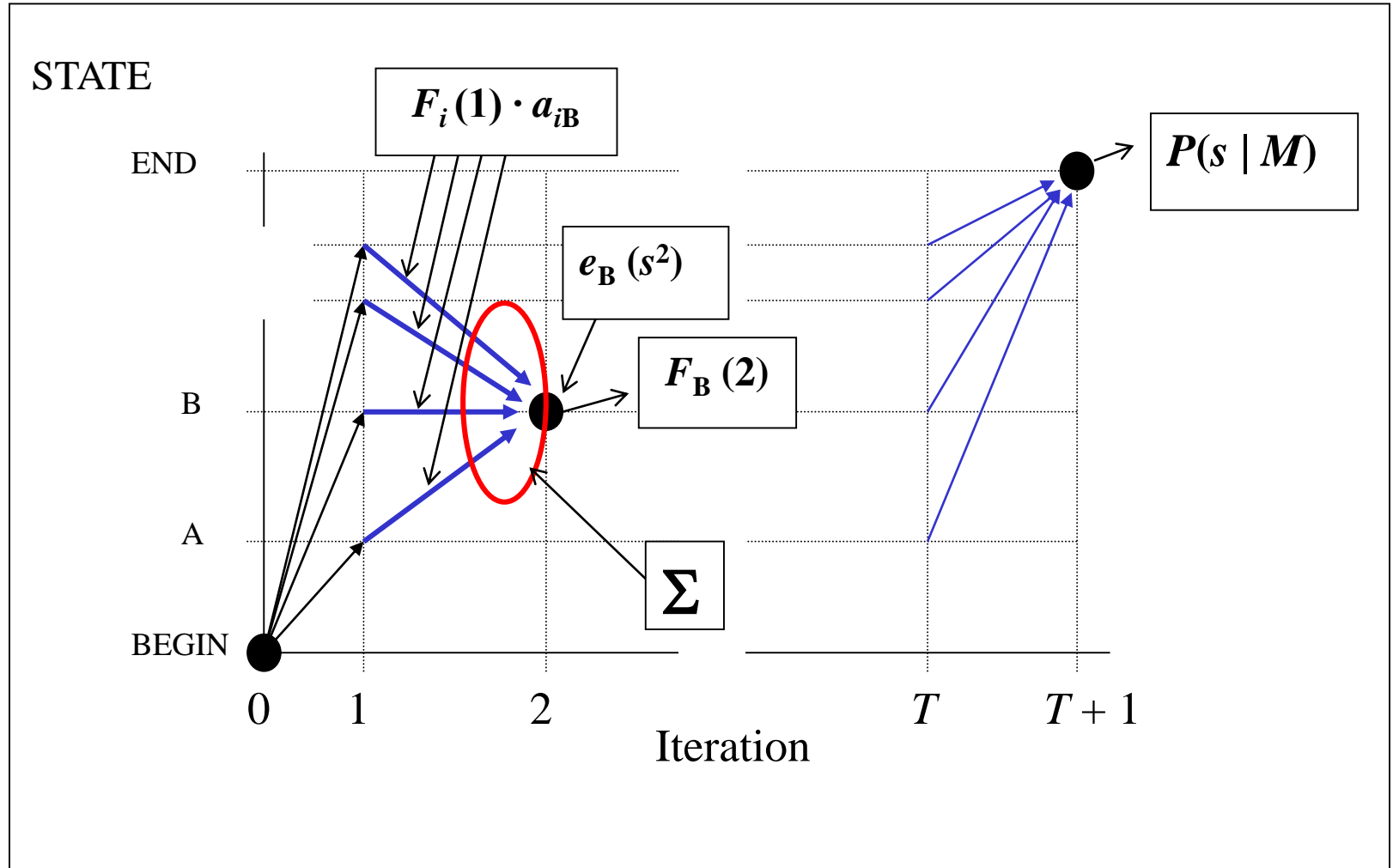
If we know the probabilities of emitting the two first characters of the sequence ending the path in states  $L$  and  $R$  respectively:

$$F_Y(2) \equiv P(s^1, s^2, \pi(2) = Y \mid M) \quad \text{and} \quad F_N(2) \equiv P(s^1, s^2, \pi(2) = N \mid M)$$

then we can compute:

$$P(s^1, s^2, s^3, \pi(3) = Y \mid M) = F_Y(2) \cdot a_{YY} \cdot e_Y(s^3) + F_N(2) \cdot a_{NY} \cdot e_Y(s^3)$$

# Forward Algorithm



## Forward algorithm: computational complexity

## *Naïf method*

$$P(s|M) = \sum_{\pi} P(s, \pi|M)$$

There are  $N^T$  possible paths.

Each path requires about  $2 \cdot T$  operations.

The time for the computation is  $\mathcal{O}(\mathbf{T} \cdot \mathbf{N}^{\mathbf{T}})$

[illegible]

Emission:  $0.1 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.1 \times 0.1 \times 0.1 \times 0.1 \times 0.4 \times 0.1 \times 0.4$

[illegible][illegible]

Emission:  $0.1 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.4 \times 0.1 \times 0.1 \times 0.1 \times 0.1 \times 0.4 \times 0.1 \times 0.25$

[illegible]

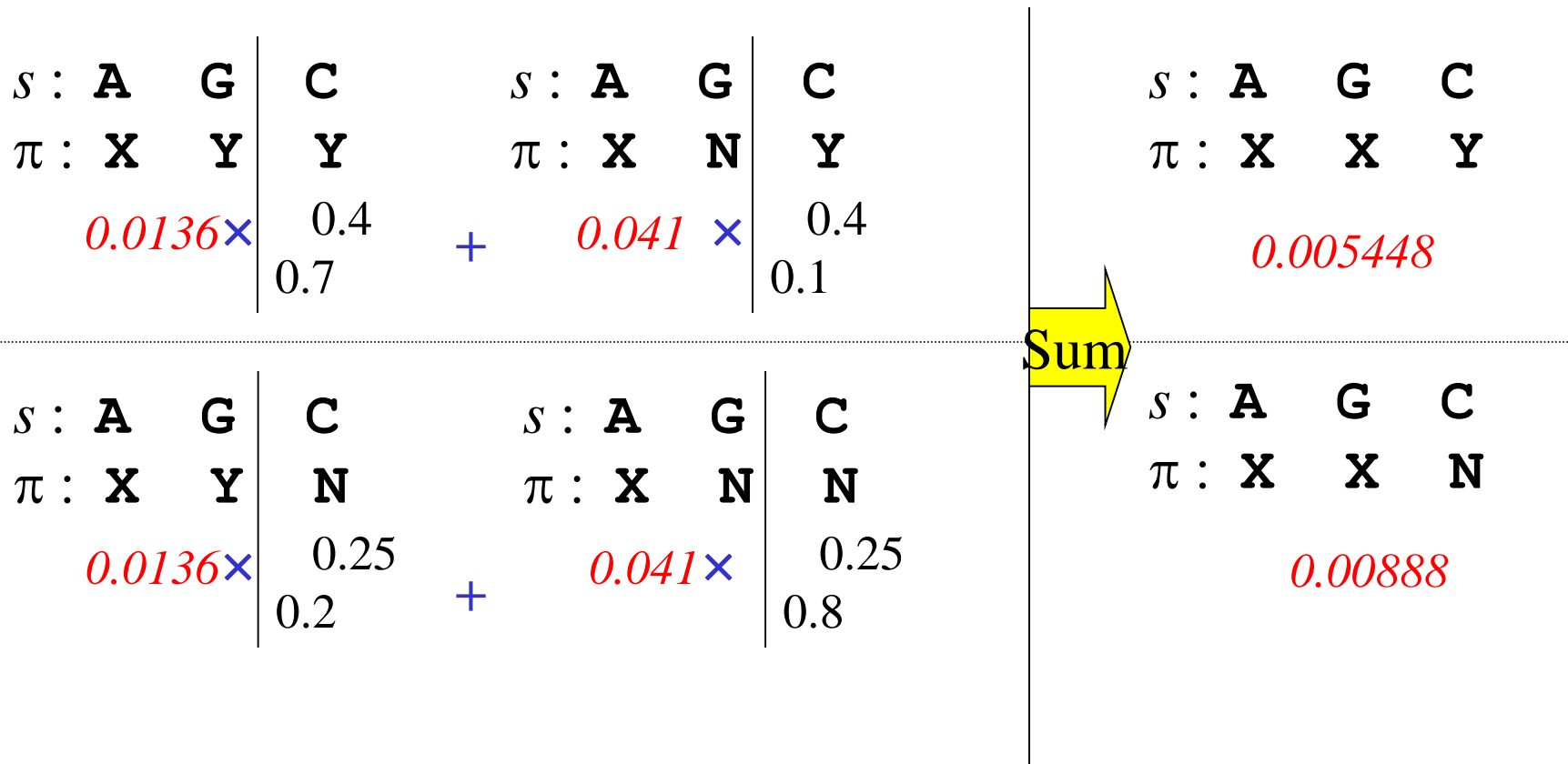
# Forward algorithm: computational complexity

- *Forward algorithm*

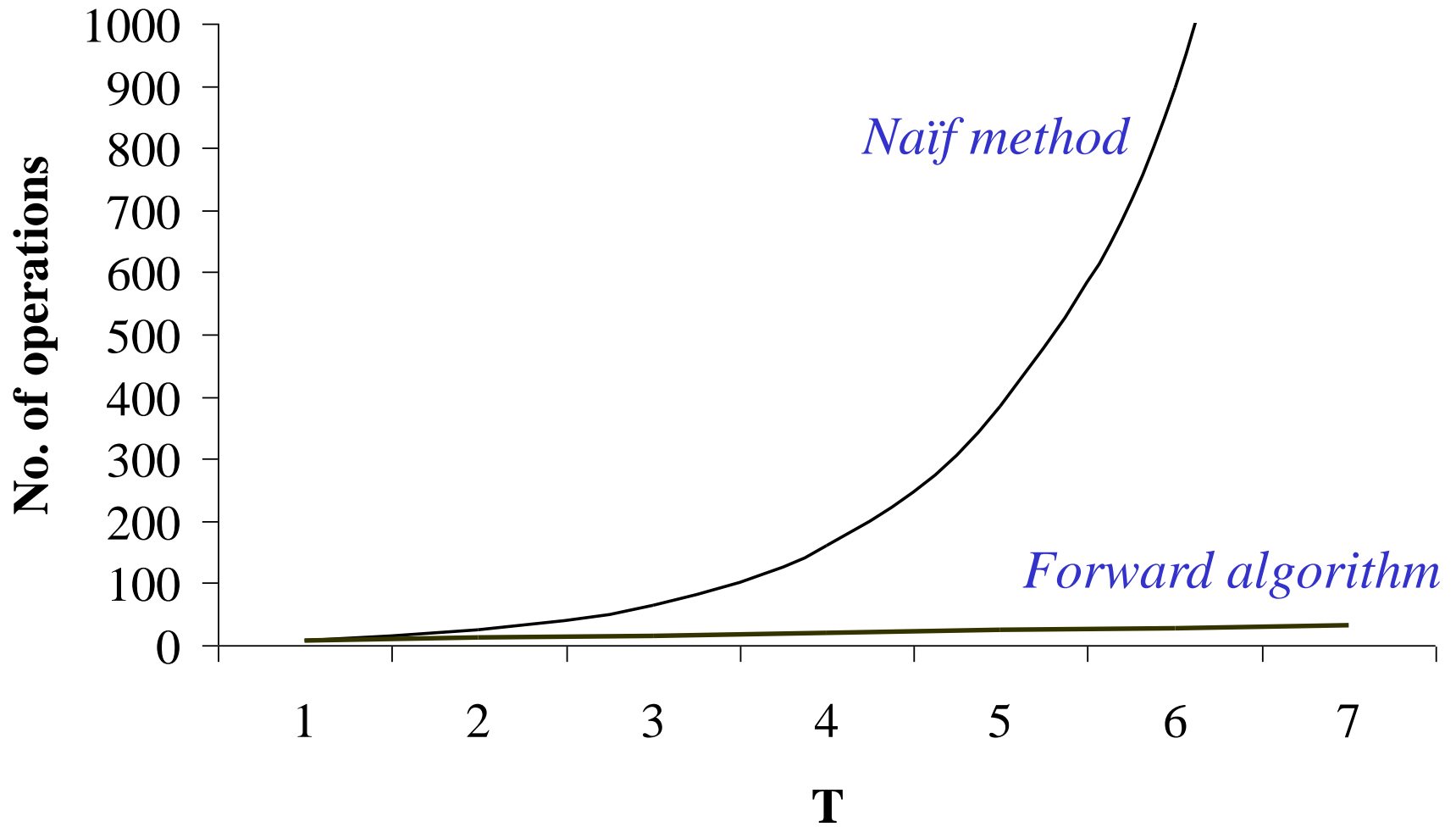
$T$  positions,  $N$  values for each position

Each element requires about  $2 \cdot N$  product and  $1$  sum

The time for the computation is  $O(T \cdot N^2)$



# Forward algorithm: computational complexity





# Hidden Markov Models: Algorithms

- Résumé
- Evaluating  $P(s \mid M)$ : Forward Algorithm
- Evaluating  $P(s \mid M)$ : Backward Algorithm

## Backward Algorithm

Similar to the Forward algorithm: it computes  $P(s / M)$ , reconstructing the sequence from the end

For each state  $k$  and each position  $i$  in the sequence, we compute:

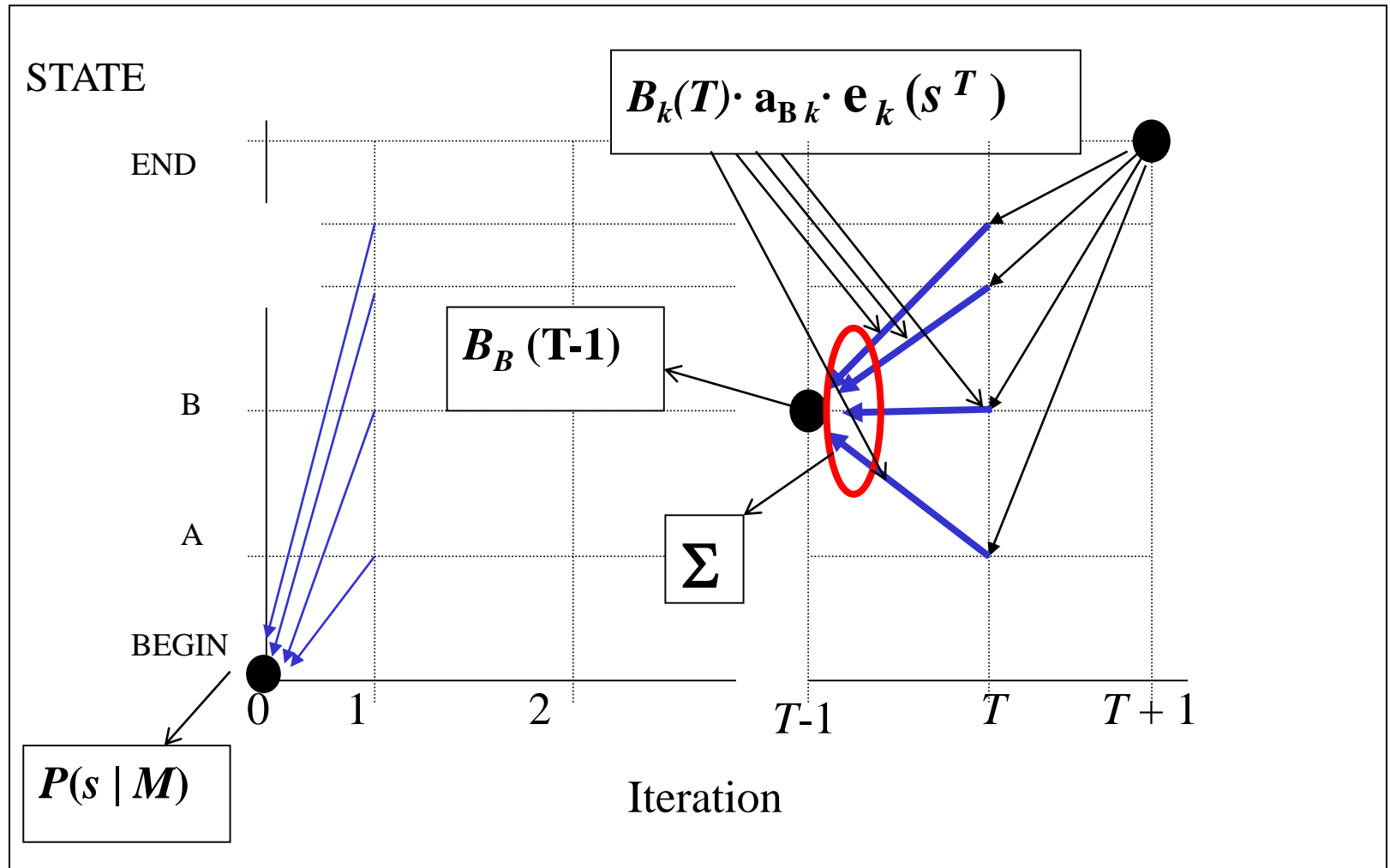
$$B_k(i) = P(s^{i+1}s^{i+2}s^{i+3} \dots s^T / \pi(i) = k)$$

*Initialisation:*  $B_k(T) = P(\pi(T+1) = \text{END} \mid \pi(T) = k) = a_{k0}$

*Recurrence:*  $B_l(i-1) = P(s^i s^{i+1} \dots s^T / \pi(i-1) = l) =$   
 $= \sum_k P(s^{i+1} s^{i+2} \dots s^T / \pi(i) = k) \cdot a_{lk} \cdot e_k(s^i) =$   
 $= \sum_k B_k(i) \cdot e_k(s^i) \cdot a_{lk}$

*Termination:*  $P(s) = P(s^1 s^2 s^3 \dots s^T / \pi(0) = \text{BEGIN}) =$   
 $= \sum_k P(s^2 \dots s^T / \pi(1) = k) \cdot a_{0k} \cdot e_k(s^1) =$   
 $= \sum_k B_k(1) \cdot a_{0k} \cdot e_k(s^1)$

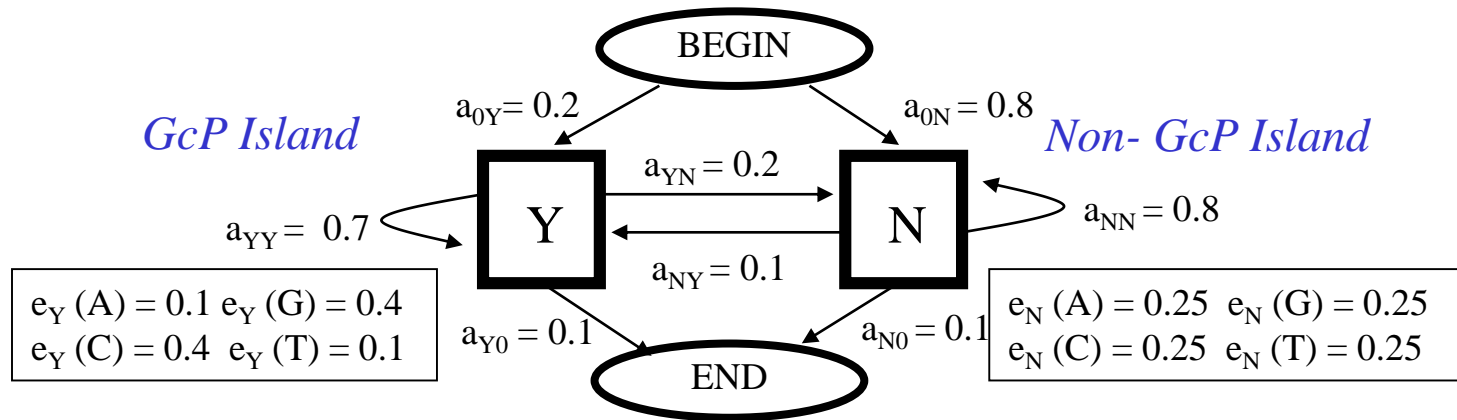
# Backward Algorithm



# Hidden Markov Models: Algorithms

- Résumé
- Evaluating  $P(s \mid M)$ : Forward Algorithm
- Evaluating  $P(s \mid M)$ : Backward Algorithm
- Showing the path: Viterbi decoding

# Finding the best path



$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{Y}$

Emission:  $0.1 \times 0.4$   
 Transition:  $0.2 \times 0.7$   
**0.0056**

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{Y}$

Emission:  $0.25 \times 0.4$   
 Transition:  $0.8 \times 0.1$   
**0.008**

**In order to continue the process we only need to know the last state:**

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{N}$

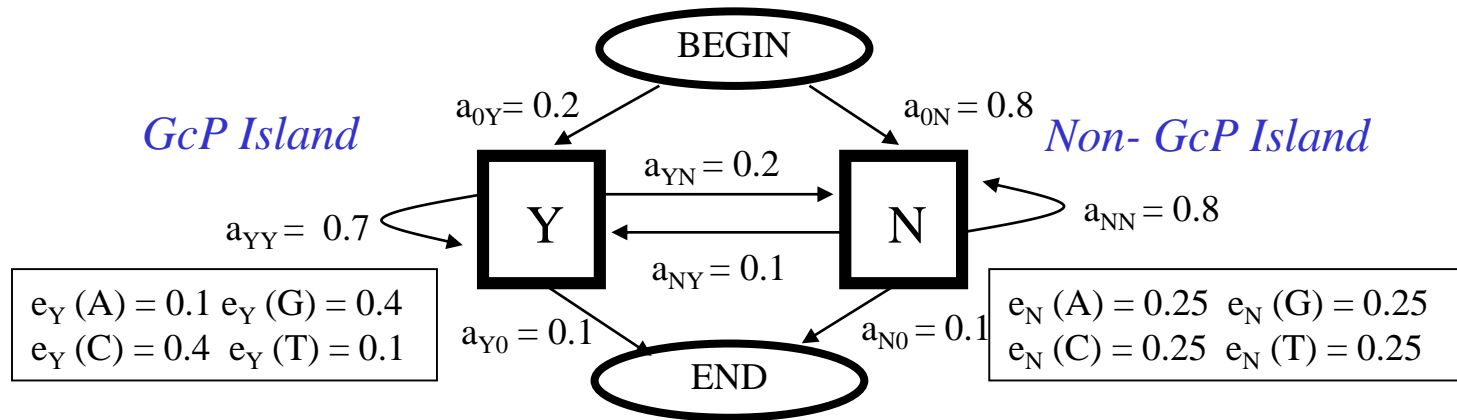
Emission:  $0.1 \times 0.25$   
 Transition:  $0.2 \times 0.2$   
**0.001**

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{N}$

Emission:  $0.25 \times 0.25$   
 Transition:  $0.8 \times 0.8$   
**0.04**

**We can choose for each final state the path scoring with the maximum probability**

# Finding the best path



$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{Y}$

Emission:  $0.1 \times 0.4$   
 Transition:  $0.2 \times 0.7$   
0.0056

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{Y}$

Emission:  $0.25 \times 0.4$   
 Transition:  $0.8 \times 0.1$   
0.008

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{Y}$

0.008

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{Y} \quad \mathbf{N}$

Emission:  $0.1 \times 0.25$   
 Transition:  $0.2 \times 0.2$   
0.001

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{N}$

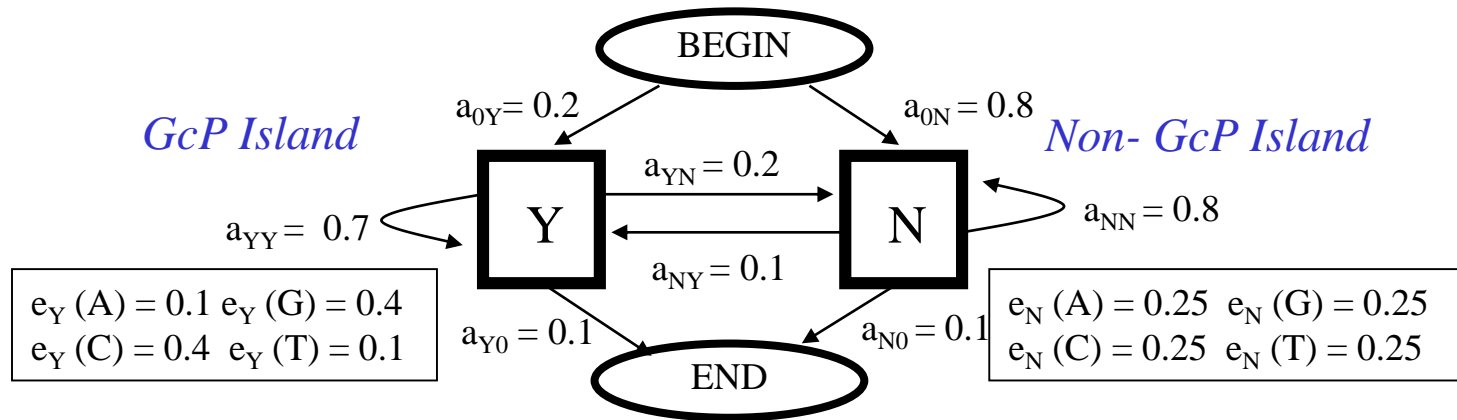
Emission:  $0.25 \times 0.25$   
 Transition:  $0.8 \times 0.8$   
0.04

Max

$s : \mathbf{A} \quad \mathbf{G}$   
 $\pi : \mathbf{N} \quad \mathbf{N}$

0.04

# Finding the best path

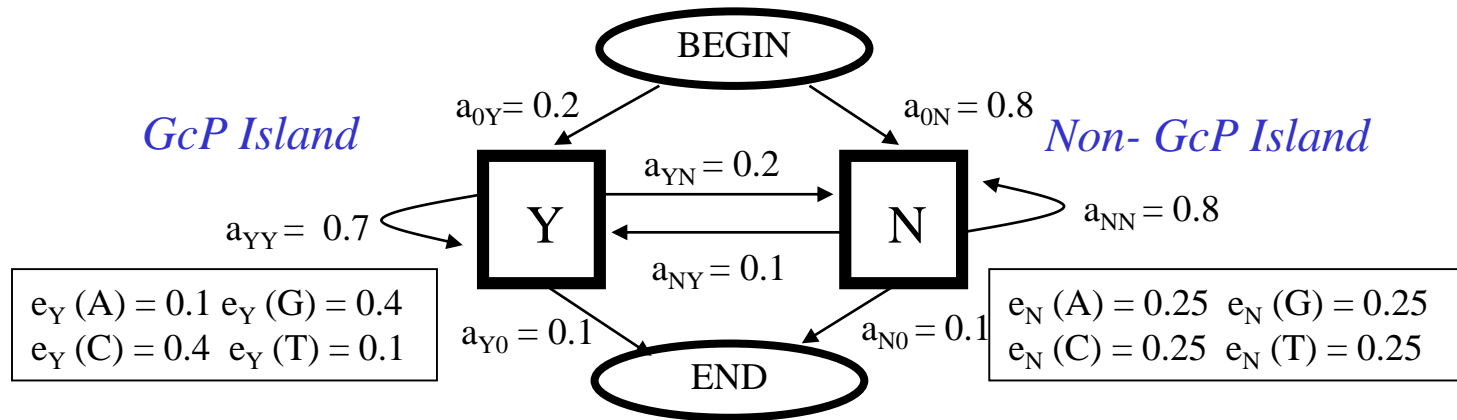


$s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C}$	$s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C}$
$\pi : \mathbf{N} \quad \mathbf{Y} \quad \mathbf{Y}$	$\pi : \mathbf{N} \quad \mathbf{N} \quad \mathbf{Y}$
$0.008 \times 0.4$	$0.04 \times 0.4$
$0.7$	$0.1$
$= \mathbf{0.00224}$	$= \mathbf{0.0016}$

$s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C}$	$s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C}$
$\pi : \mathbf{N} \quad \mathbf{Y} \quad \mathbf{N}$	$\pi : \mathbf{N} \quad \mathbf{N} \quad \mathbf{N}$
$0.008 \times 0.25$	$0.04 \times 0.25$
$0.2$	$0.8$
$= \mathbf{0.0004}$	$= \mathbf{0.008}$

;

# Finding the best path



$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{Y} \quad \mathbf{Y} \\
 0.008 \times 0.4 \\
 \quad \quad 0.7 \\
 = \mathbf{0.00224}
 \end{array}$$

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{N} \quad \mathbf{Y} \\
 0.04 \times 0.4 \\
 \quad \quad 0.1 \\
 = \mathbf{0.0016}
 \end{array}$$

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{Y} \quad \mathbf{Y} \\
 \mathbf{0.00224}
 \end{array}$$

$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{Y} \quad \mathbf{N} \\
 0.008 \times 0.25 \\
 \quad \quad 0.2 \\
 = \mathbf{0.0004}
 \end{array}$$

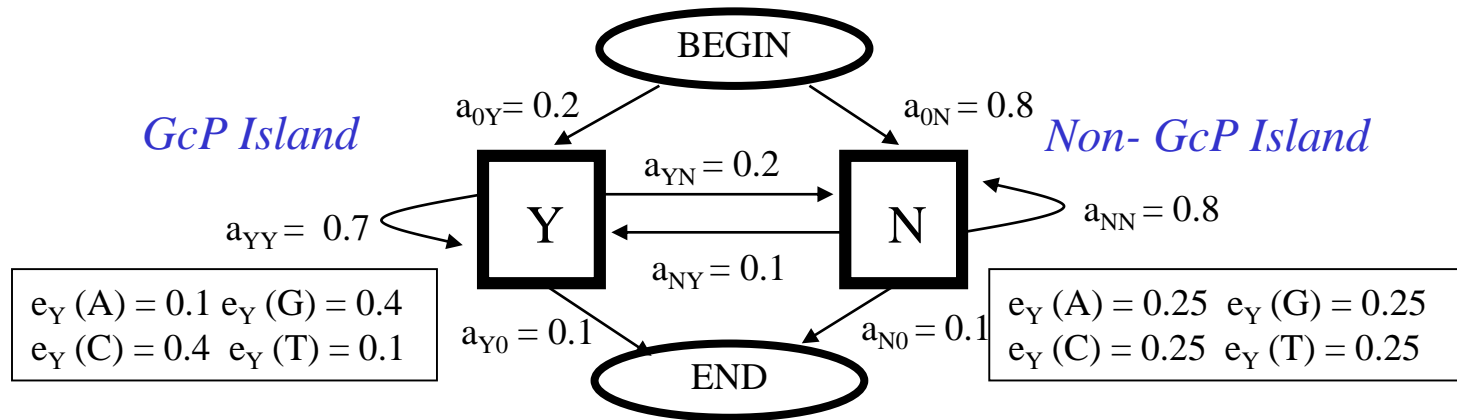
$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{N} \quad \mathbf{N} \\
 0.04 \times 0.25 \\
 \quad \quad 0.8 \\
 = \mathbf{0.008}
 \end{array}$$



$$\begin{array}{c}
 s : \mathbf{A} \quad \mathbf{G} \quad \mathbf{C} \\
 \pi : \mathbf{N} \quad \mathbf{N} \quad \mathbf{N} \\
 \mathbf{0.008}
 \end{array}$$



# Finding the best path



Iterating until the last position of the sequence:

A	G	C	G	C	G	T	A	A	T	C	T	G	
N	Y	Y	Y	Y	Y	Y	N	N	N	Y	Y	Y	0.1 ( $a_{Y0}$ )
													×
A	G	C	G	C	G	T	A	A	T	C	T	G	
N	N	N	N	N	N	N	N	N	N	N	N	N	0.1 ( $a_{N0}$ )
													×

Choose the Maximum

# Viterbi Algorithm

$$\pi^* = \operatorname{argmax}_{\pi} [ P( \pi, s \mid M ) ]$$

The computation of  $P(s, \pi^* \mid M)$  can be decomposed in simplest problems

Let  $V_k(i)$  be the probability of the most probable path for generating the subsequence  $s^1 s^2 s^3 \dots s^i$  ending in the state  $k$  at iteration  $i$

*Initialisation:*  $V_{BEGIN}(0) = 1 \quad V_i(0) = 0 \quad \forall \ i \neq \text{BEGIN}$

*Recurrence:*  $V_l(i+1) = e_l(s^{i+1}) \cdot \operatorname{Max}_k ( V_k(i) \cdot a_{kl} )$

$$\operatorname{ptr}_i(l) = \operatorname{argmax}_k ( V_k(i) \cdot a_{kl} )$$

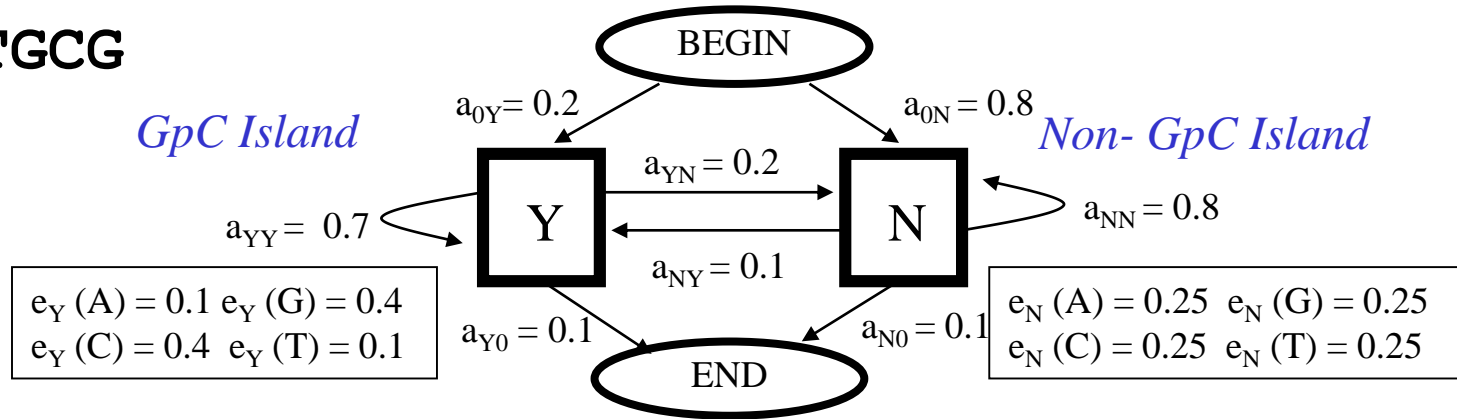
*Termination:*  $P(s, \pi^*) = \operatorname{Max}_k ( V_k(T) \cdot a_{k0} )$

$$\pi^*(T) = \operatorname{argmax}_k ( V_k(T) \cdot a_{k0} )$$

*Traceback:*  $\pi^*(i-1) = \operatorname{ptr}_i(\pi^*(i))$

# Viterbi Algorithm: Example

**S = ATGCG**

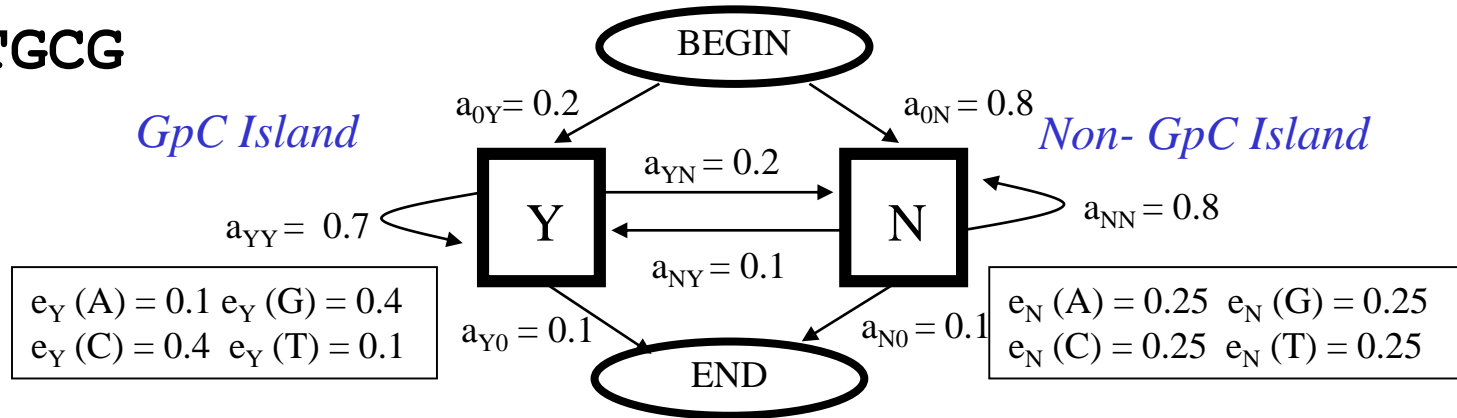


*Initialisation:*  $V_{BEGIN}(0) = 1$   $V_i(0) = 0 \quad \forall i \neq BEGIN$

	-	A	T	G	C	G	-
Begin	1						
Y	0						
N	0						
End	0						

# Viterbi Algorithm: Example

**S = ATGCG**



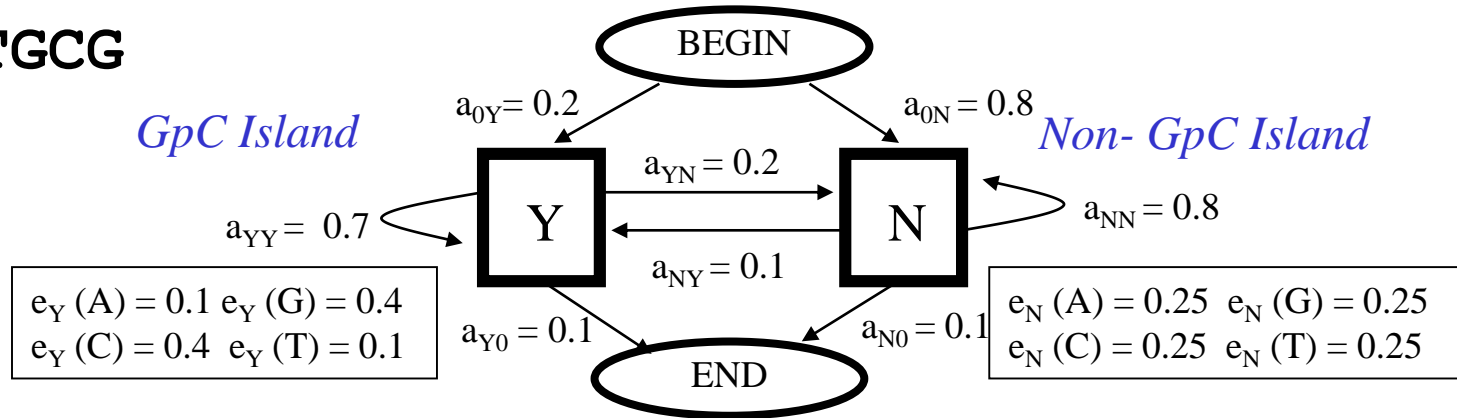
*Recurrence:*  $V_l(1) = e_l(s^1) \cdot \text{Max}_k (V_k(0) \cdot a_{kl})$

$\text{ptr}_1(l) = \text{argmax}_k (V_k(0) \cdot a_{kl})$

	-	A	T	G	C	G	-
Begin	1	0					
Y	0	0.2x0.1= =2e-2 ptr=Begin					
N	0	0.8x0.25= =0.2 ptr=Begin					
End	0	0					

# Viterbi Algorithm: Example

**S = ATGCG**



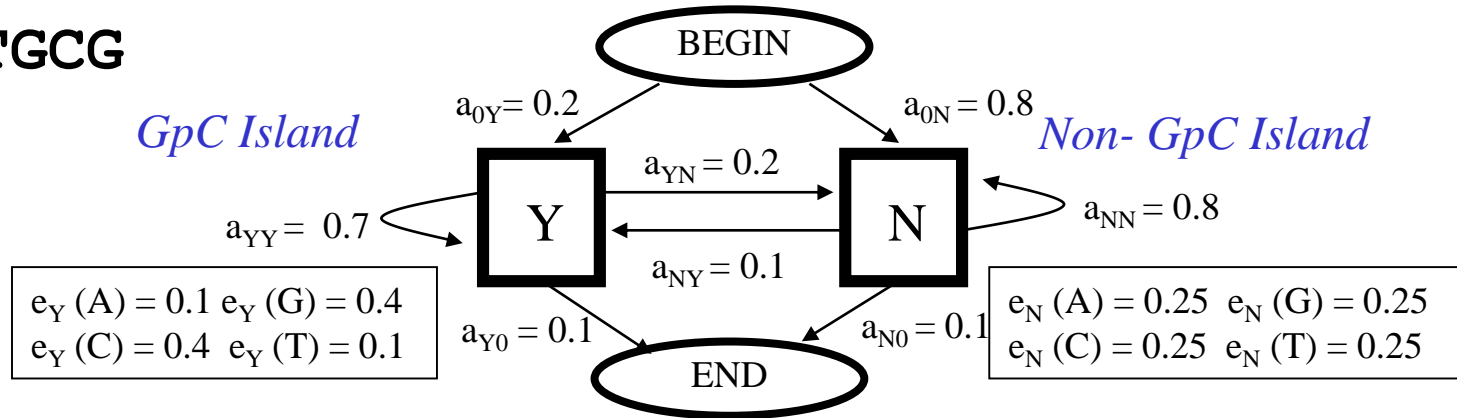
*Recurrence:*  $V_l(2) = e_l(s^2) \cdot \text{Max}_k (V_k(1) \cdot a_{kl})$

$\text{ptr}_2(l) = \text{argmax}_k (V_k(1) \cdot a_{kl})$

	-	A	T	G	C	G	-
Begin	1	0	0				
Y	0	2e-2 Begin	Max(2e-2x0.7x0.1; <b>0.2x0.1x0.1</b> ) <b>2e-3; ptr=N</b>				
N	0	0.2 Begin	Max(2e-2x0.2x0.25; <b>0.2x0.8x0.25</b> ) <b>1.6e-2; ptr=N</b>				
End	0	0	0				

# Viterbi Algorithm: Example

**S = ATGCG**



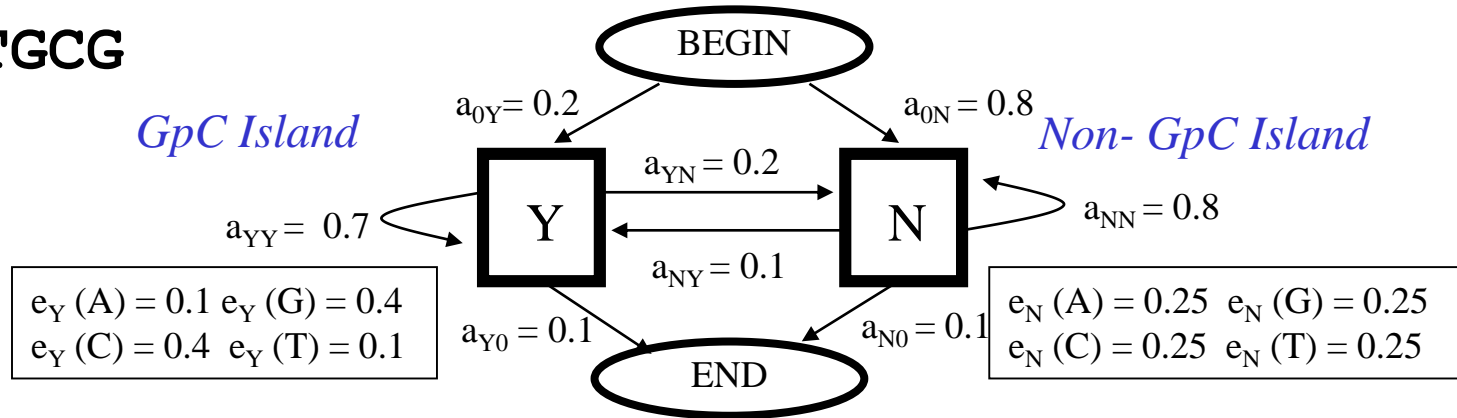
*Recurrence:*  $V_l(3) = e_l(s^3) \cdot \text{Max}_k (V_k(2) \cdot a_{kl})$

$\text{ptr}_3(l) = \text{argmax}_k (V_k(2) \cdot a_{kl})$

	-	A	T	G	C	G	-
<b>Begin</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>0</b>			
<b>Y</b>	<b>0</b>	<b>2e-2</b> <b>Begin</b>	<b>2e-3;</b> <b>ptr=N</b>	<b>Max(2e-3x0.7x0.4;</b> <b>1.6e-2x0.1x0.4)</b> <b>6.4e-4; ptr=N</b>			
<b>N</b>	<b>0</b>	<b>0.2</b> <b>Begin</b>	<b>1.6e-2;</b> <b>ptr=N</b>	<b>Max(2e-3x0.2x0.25;</b> <b>1.6e-2x0.8x0.25)</b> <b>3.2e-4; ptr=N</b>			
<b>End</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>			

# Viterbi Algorithm: Example

**S = ATGCG**



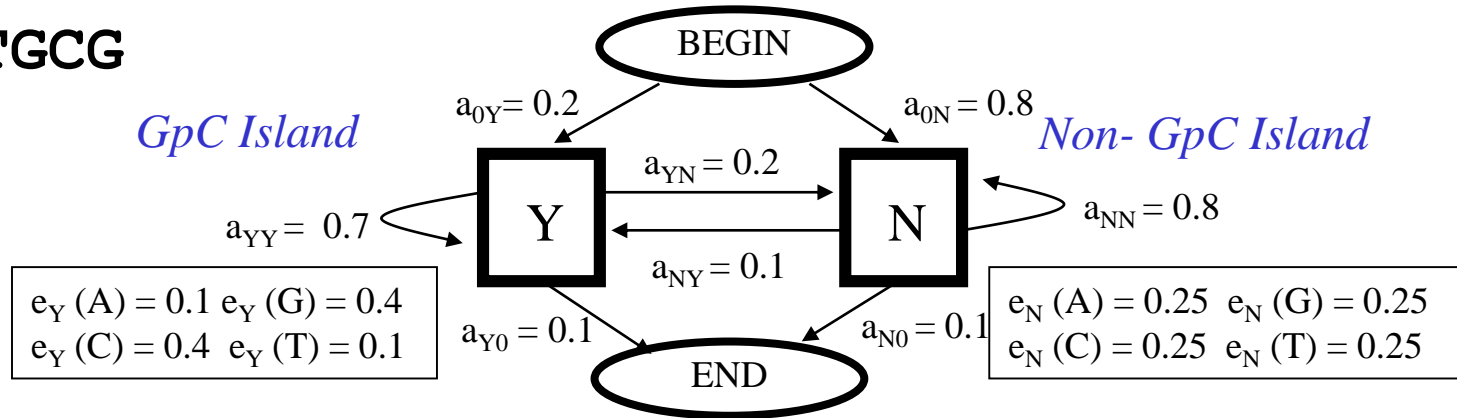
*Recurrence:*  $V_l(4) = e_l(s^4) \cdot \text{Max}_k (V_k(3) \cdot a_{kl})$

$\text{ptr}_4(l) = \text{argmax}_k (V_k(3) \cdot a_{kl})$

	-	A	T	G	C	G	-
<b>Begin</b>	1	0	0	0	0		
<b>Y</b>	0	2e-2 Begin	2e-3; ptr=N	6.4e-4; ptr=N	Max(6.4e-4x0.7x0.4; 3.2e-4x0.1x0.4) 1.792e-4; ptr=Y		
<b>N</b>	0	0.2 Begin	1.6e-2; ptr=N	3.2e-4; ptr=N	Max(6.4e-4x0.2x0.25; 3.2e-4x0.8x0.25) 6.4e-5; ptr=N		
<b>End</b>	0	0	0	0	0		

# Viterbi Algorithm: Example

**S = ATGCG**



*Recurrence:*  $V_l(5) = e_l(s^5) \cdot \text{Max}_k (V_k(4) \cdot a_{kl})$

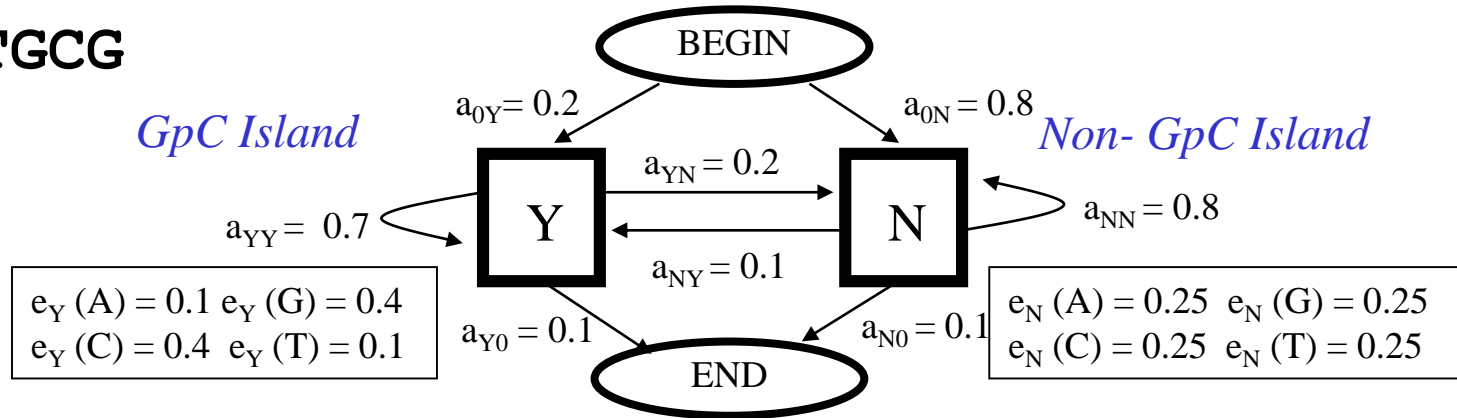
$\text{ptr}_5(l) = \text{argmax}_k (V_k(4) \cdot a_{kl})$

	-	A	T	G	C	G	-
Begin	1	0	0	0	0	0	
Y	0	2e-2 Begin	2e-3; ptr=N	6.4e-4; ptr=N	1.792e-4; ptr=Y	Max(1.792e-4x0.7x0.4; 6.4e-5x0.1x0.4) 5.0176e-5; ptr=Y	
N	0	0.2 Begin	1.6e-2; ptr=N	3.2e-4; ptr=N	6.4e-5; ptr=N	Max(1.792e-4x0.2x0.25; 6.4e-5x0.8x0.25) 1.28e-5; ptr=N	
End	0	0	0	0	0	0	



# Viterbi Algorithm: Example

**S = ATGCG**



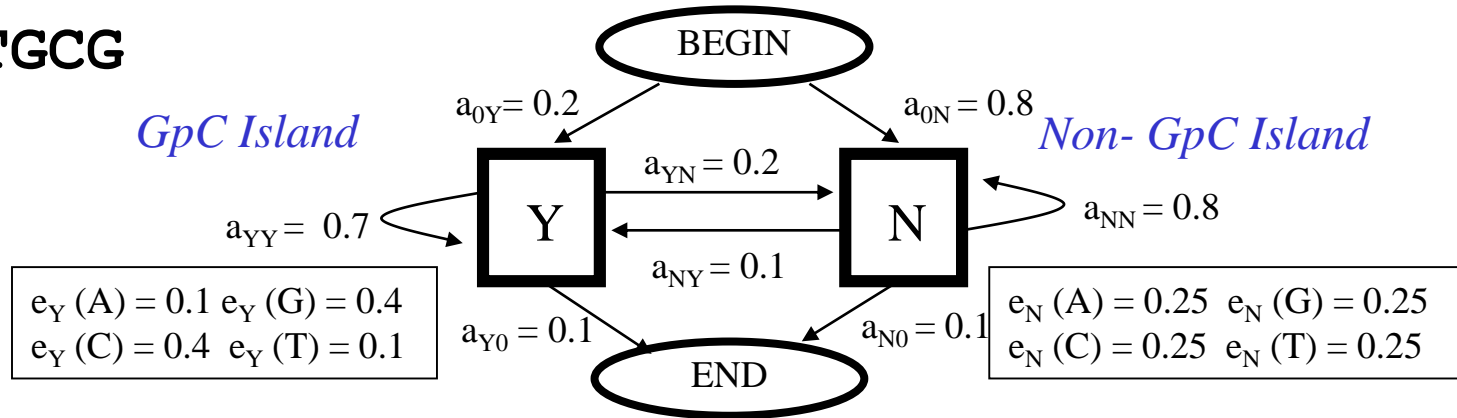
*Termination:*  $P(s, \pi^*) = \text{Max}_k (V_k(T) \cdot a_{k0})$

$$\pi^*(T) = \text{argmax}_k (V_k(T) \cdot a_{k0})$$

	-	A	T	G	C	G	-
<b>Begin</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
<b>Y</b>	<b>0</b>	<b>2e-2</b> <b>Begin</b>	<b>2e-3;</b> <b>ptr=N</b>	<b>6.4e-4;</b> <b>ptr=N</b>	<b>1.792e-4;</b> <b>ptr=Y</b>	<b>5.0176e-5;</b> <b>ptr=Y</b>	<b>0</b>
<b>N</b>	<b>0</b>	<b>0.2</b> <b>Begin</b>	<b>1.6e-2;</b> <b>ptr=N</b>	<b>3.2e-4;</b> <b>ptr=N</b>	<b>6.4e-5;</b> <b>ptr=N</b>	<b>1.28e-5;</b> <b>ptr=N</b>	<b>0</b>
<b>End</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>Max(5.0176e-5x0.1;</b> <b>1.28e-5x0.1)</b> <b>5.0176e-6; ptr=Y</b>

# Viterbi Algorithm: Example

**S = ATGCG**



*Traceback:*  $\pi^* (i-1) = \text{ptr}_i (\pi^* (i))$

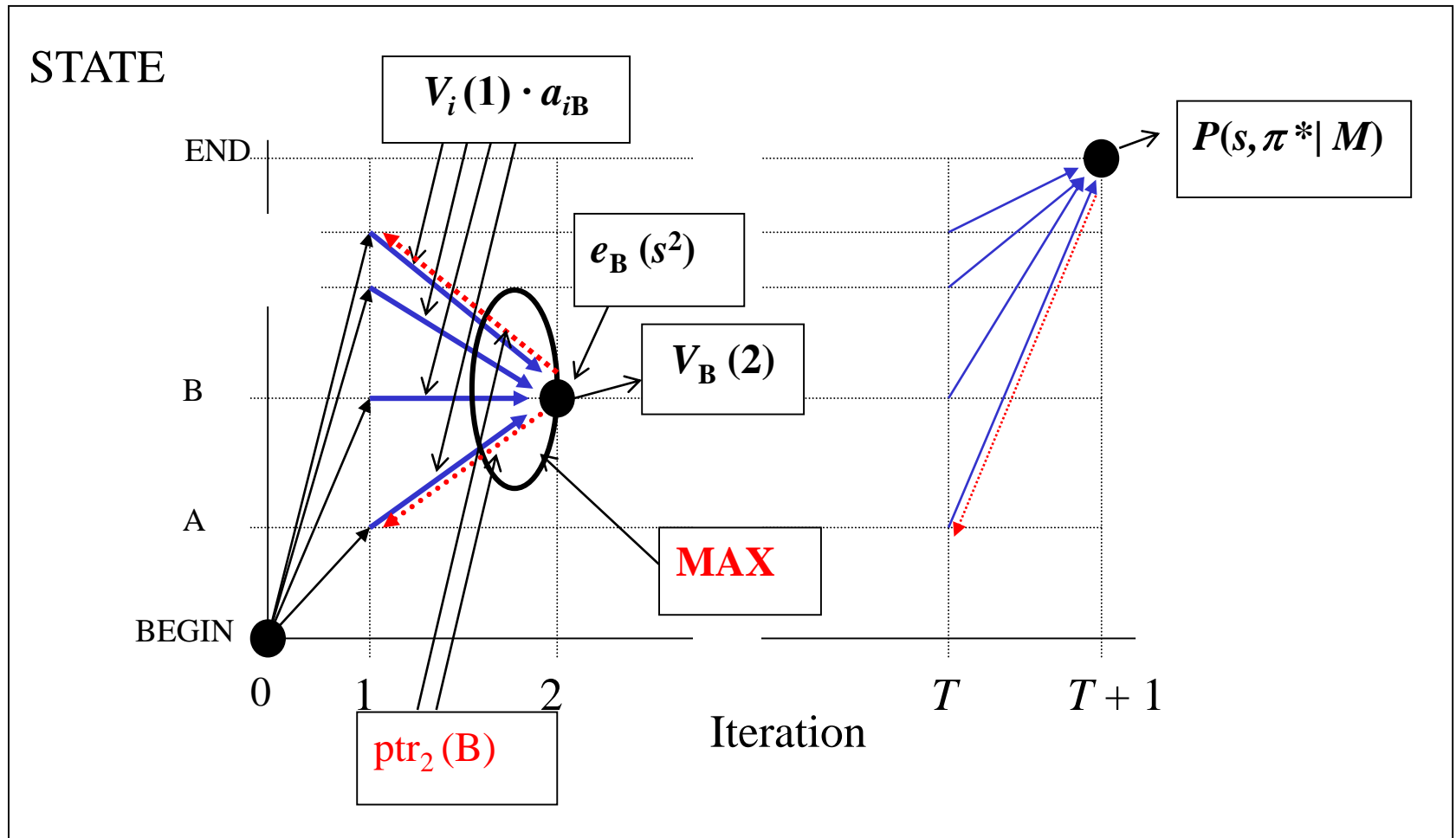
	-	A	T	G	C	G	-
Begin	1	0	0	0	0	0	0
Y	0	2e-2 Begin	2e-3; ptr=N	6.4e-4; ptr=N	1.792e-4; ptr=Y	5.0176e-5; ptr=Y	0
N	0	0.2 Begin	1.6e-2; ptr=N	3.2e-4; ptr=N	6.4e-5; ptr=N	1.28e-5; ptr=N	0
End	0	0	0	0	0	0	5.0176e-6; ptr=Y

**S = - A T G C G -**

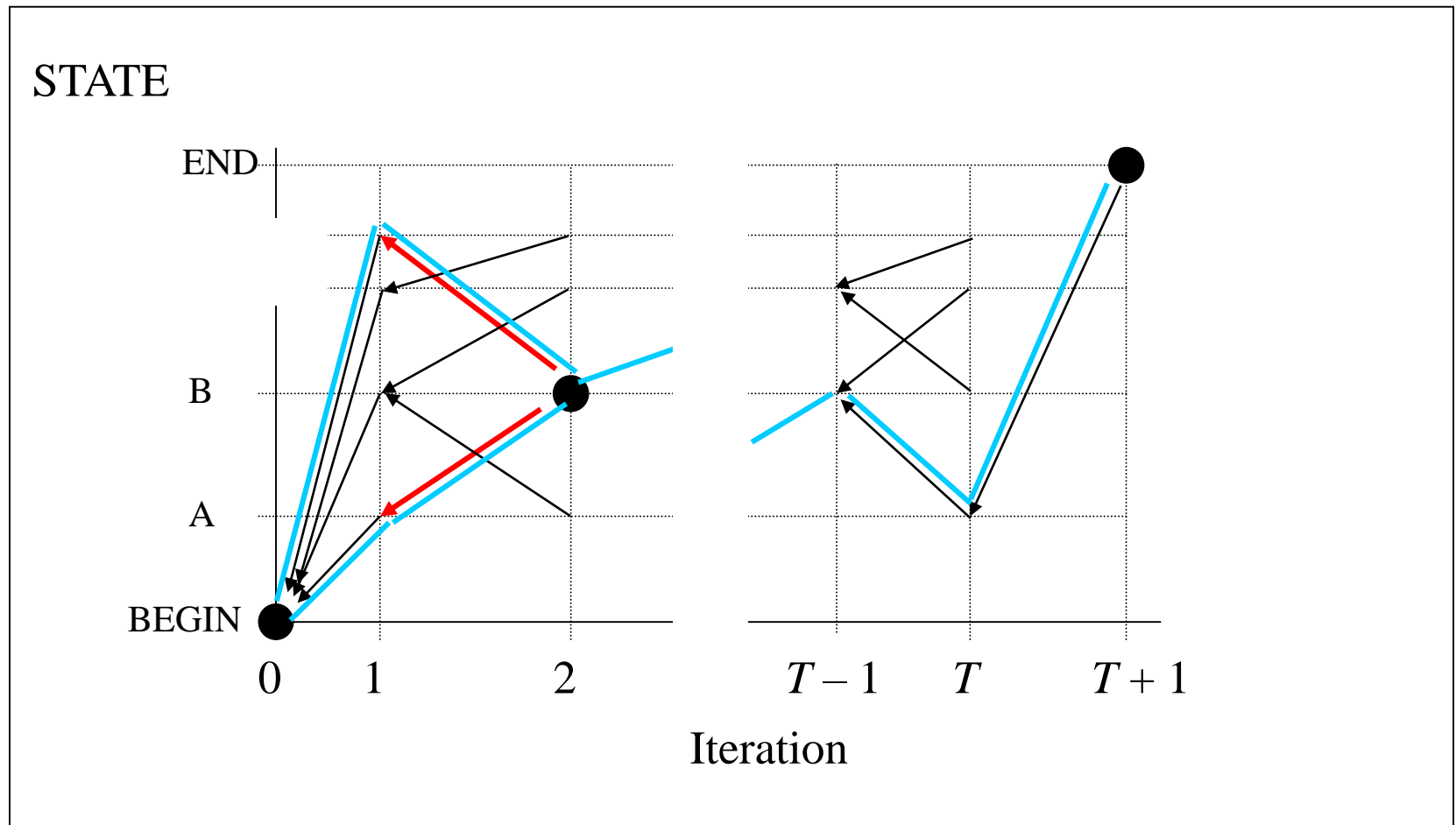
**$\pi^* = B N N Y Y Y E$**

$P(s, \pi^* | M) = 5.0176e-6$

# Viterbi Algorithm



# Viterbi Algorithm



Viterbi path —————

Different paths can have the same probability

## An alternative: A posteriori decoding

For each position choose the state  $\underline{\pi}(t)$  :

$$\underline{\pi}(i) = \operatorname{argmax}_k [ P( \pi(i) = k | s, M ) ]$$

How to compute  $P( \pi(i) = k | s, M )$  for any state  $k$  and any position  $i$ ?

$$P(\pi(i) = k | s, M) = \frac{P(\pi(i) = k, s | M)}{P(s | M)}$$

$$\begin{aligned} P(\pi(i) = k, s | M) &= P(s^1 s^2 \dots s^i, \pi(i) = k | M) \cdot P(s^{i+1}, s^{i+2}, \dots s^T | \pi(i) = k, M) = \\ &= F_k(i) \cdot B_k(i) \end{aligned}$$

$$P(\pi(i) = k | s, M) = \frac{F_k(i) \cdot B_k(i)}{P(s | M)}$$

Elements of the Forward and Backward matrices

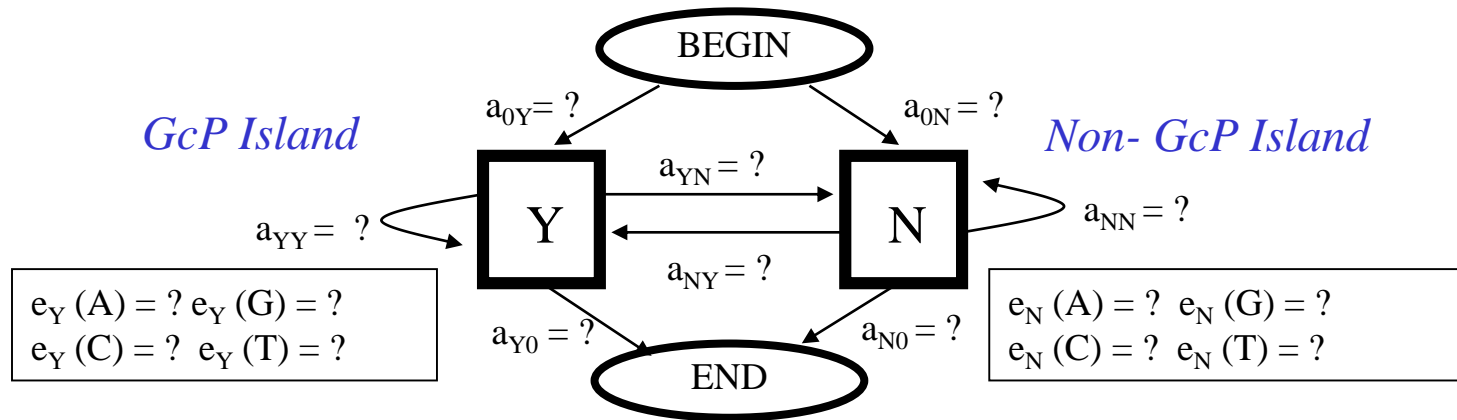
Computed with Forward or Backward algorithm termination steps

**NB: this decoding can contain “non allowed” transitions**

# Hidden Markov Models: Algorithms

- Résumé
- Evaluating  $P(s \mid M)$ : Forward Algorithm
- Evaluating  $P(s \mid M)$ : Backward Algorithm
- Showing the path: Viterbi decoding
- Showing the path: A posteriori decoding
- Training a model: EM algorithm

# If we know the path generating the training sequence



$S : \mathbf{A \quad G \quad C \quad G \quad C \quad G \quad T \quad A \quad A \quad T \quad C \quad T \quad G}$

$\pi : \mathbf{Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad N \quad N \quad N \quad N \quad N \quad N}$

Emission:  $e_Y(A) \times e_Y(G) \times e_Y(A) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(T) \times e_N(A) \times e_N(A) \times e_N(T) \times e_N(C) \times e_N(T) \times e_N(G)$

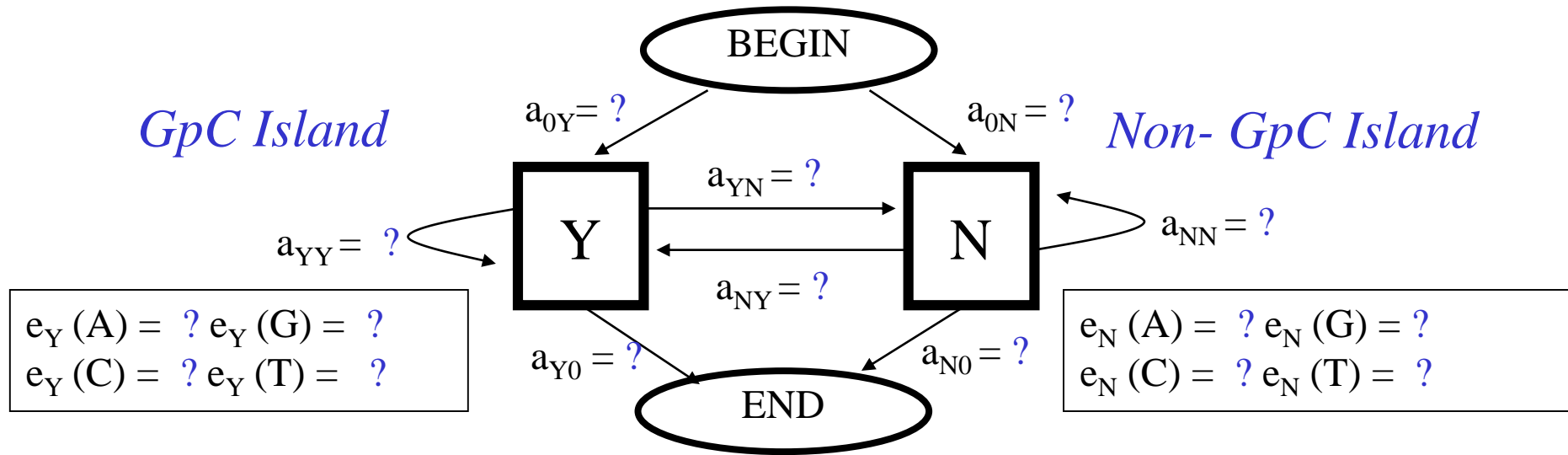
Transition:  $a_{0Y} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{N0}$

Just count!

Example:  $a_{YY} = n_{YY} / (n_{YY} + n_{YN}) = 6/7$

$e_Y(A) = n_Y(A) / [n_Y(A) + n_Y(C) + n_Y(G) + n_Y(T)] = 1/7$

## GpC Island, simple model



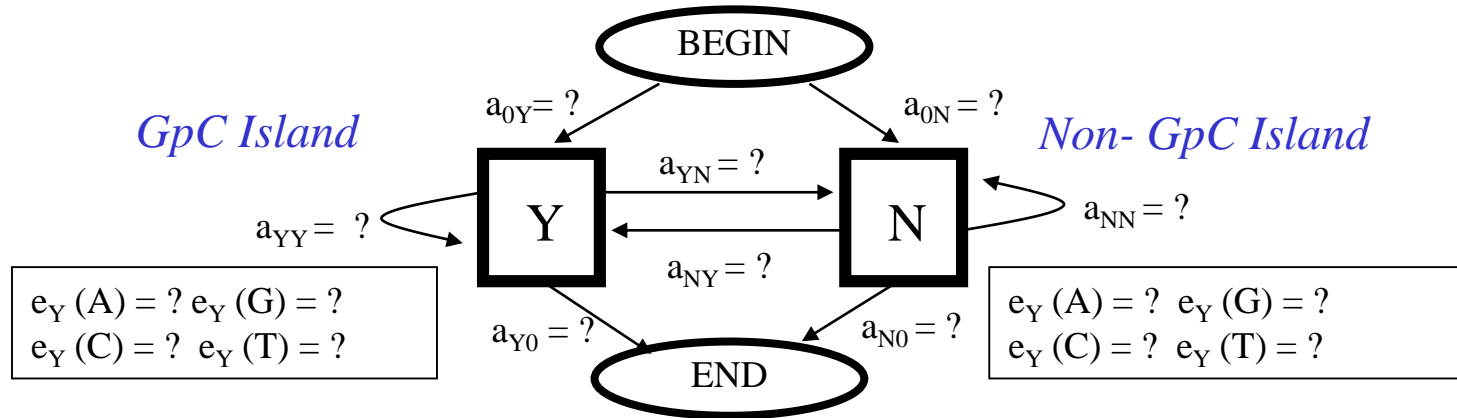
$s$  : **AGCGCGTAATCTG**

$\pi$  : **YYYYYYNNNNNN**

- $P(s, \pi / M)$  can be easily computed
- How to evaluate  $P(s / M)$ , when the path is unknown?
- Can we show the hidden path?
- Can we evaluate the parameters starting from known examples?



# Can we evaluate the parameters starting from known examples?



$S : \mathbf{A \quad G \quad C \quad G \quad C \quad G \quad T \quad A \quad A \quad T \quad C \quad T \quad G}$

$\pi : \mathbf{Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad N \quad N \quad N \quad N \quad N \quad N}$

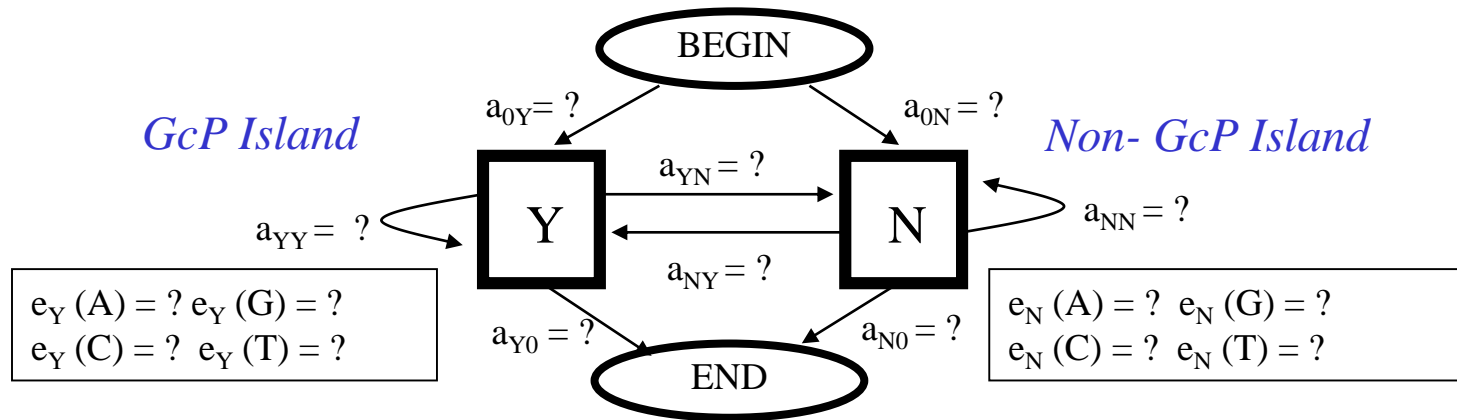
Emission:  $e_Y(A) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(T) \times e_N(A) \times e_N(A) \times e_N(T) \times e_N(C) \times e_N(T) \times e_N(G)$

Transition:  $a_{0Y} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{N0}$

How to find the parameters  $e$  and  $a$  that maximises this probability?

How if we don't know the path?

# If we know the path generating the training sequence



$S : \mathbf{A \quad G \quad C \quad G \quad C \quad G \quad T \quad A \quad A \quad T \quad C \quad T \quad G}$

$\pi : \mathbf{Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad Y \quad N \quad N \quad N \quad N \quad N \quad N}$

Emission:  $e_Y(A) \times e_Y(G) \times e_Y(A) \times e_Y(G) \times e_Y(C) \times e_Y(G) \times e_Y(T) \times e_N(A) \times e_N(A) \times e_N(T) \times e_N(C) \times e_N(T) \times e_N(G)$

Transition:  $a_{0Y} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YY} \times a_{YN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{NN} \times a_{N0}$

Just count!

Example:  $a_{YY} = n_{YY} / (n_{YY} + n_{YN}) = 6/7$

$e_Y(A) = n_Y(A) / [n_Y(A) + n_Y(C) + n_Y(G) + n_Y(T)] = 1/7$



## Baum-Welch algorithm (simple discussion)

$s :$	<b>A</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>G</b>	<b>T</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>C</b>	<b>T</b>	<b>G</b>
$\pi_1 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>
$\pi_2 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>
$\pi_3 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>Y</b>
$\pi_4 :$	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>Y</b>	<b>N</b>	<b>N</b>
.....													

Given a path  $\pi$  we can count

the number of transition between states  $k$  and  $l$ :  $A_{k,l}(\pi)$

the number on emissions of character  $c$  from state  $k$ :  $E_k(c, \pi)$

We can compute the expected values over all the paths, given initial parameters  $\theta^0$

$$A_{k,l} = \sum_{\pi} P(\pi \mid s, \theta^0) \cdot A_{k,l}(\pi)$$

$$E_k(c) = \sum_{\pi} P(\pi \mid s, \theta^0) \cdot E_k(c, \pi)$$

The updated parameters are:

$$a_{k,l} = \frac{A_{k,l}}{\sum_{m=1}^N A_{k,m}}$$

$$e_k(c) = \frac{E_k(c)}{\sum_c E_k(c)}$$

Then we can iterate...

## Baum-Welch implementation

How to compute the expected number of transitions and emissions over all the paths

$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

$$B_k(i) = P(s^{i+1} s^{i+2} s^{i+3} \dots s^T / \pi(i) = k)$$

$$A_{k,l} = \sum_i P(\pi(i) = k, \pi(i+1) = l / s, \theta) = \frac{\sum_i F_k(i) \cdot a_{kl} \cdot e_l(s^{i+1}) \cdot B_l(i+1)}{P(s)}$$

$$E_k(c) = \sum_i P(s^i = c, \pi(i) = k / s, \theta) = \frac{\sum_{s^i = c} F_k(i) \cdot B_k(i)}{P(s)}$$

## WHY? Expectation-Maximisation algorithm

We need to estimate the Maximum Likelihood parameters when the paths generating the training sequences are unknown

$$\theta^{\text{ML}} = \operatorname{argmax}_{\theta} [\text{P} ( s \mid \theta, M )]$$

Given a model with parameters  $\theta^0$  the EM algorithm finds new parameters  $\theta$  that increase the likelihood of the model:

$$\text{P}( s \mid \theta ) > \text{P}( s \mid \theta^0 )$$

## WHY? Expectation-Maximisation algorithm

We need to estimate the Maximum Likelihood parameters when the paths generating the training sequences are unknown

$$\theta^{\text{ML}} = \operatorname{argmax}_{\theta} [\text{P} ( s \mid \theta, M )]$$

Given a model with parameters  $\theta^0$  the EM algorithm finds new parameters  $\theta$  that increase the likelihood of the model:

$$\text{P}( s \mid \theta ) > \text{P}( s \mid \theta^0 )$$

or equivalently

$$\ln \text{P}( s \mid \theta ) > \ln \text{P}( s \mid \theta^0 )$$

## WHY? Expectation-Maximisation algorithm

$$\ln P(s | \theta) = \ln P(s, \pi | \theta) - \ln P(\pi | s, \theta)$$

Multiplying for  $P(\pi | s, \theta^0)$  and summing over all the possible paths:

$$\ln P(s | \theta) = \sum_{\pi} P(\pi | s, \theta^0) \cdot \ln P(s, \pi | \theta) - \sum_{\pi} P(\pi | s, \theta^0) \cdot \ln P(\pi | s, \theta)$$



$Q(\theta | \theta^0)$  : Expectation value of  $\log P(s, \pi | \theta)$  over all the “current” paths

$$\begin{aligned} \ln P(s | \theta) - \ln P(s | \theta^0) &= \\ &= Q(\theta | \theta^0) - Q(\theta^0 | \theta^0) + \sum_{\pi} -P(\pi | s, \theta^0) \cdot \ln \frac{P(\pi | s, \theta)}{P(\pi | s, \theta^0)} \end{aligned}$$



# WHY? Expectation-Maximisation algorithm

$$\begin{aligned} \ln P(s | \theta) - \ln P(s | \theta^0) &= \\ &= Q(\theta | \theta^0) - Q(\theta^0 | \theta^0) + \sum_{\pi} -P(\pi | s, \theta^0) \cdot \ln \frac{P(\pi | s, \theta)}{P(\pi | s, \theta^0)} \end{aligned}$$

$\geq 0$

$$\geq Q(\theta | \theta^0) - Q(\theta^0 | \theta^0)$$

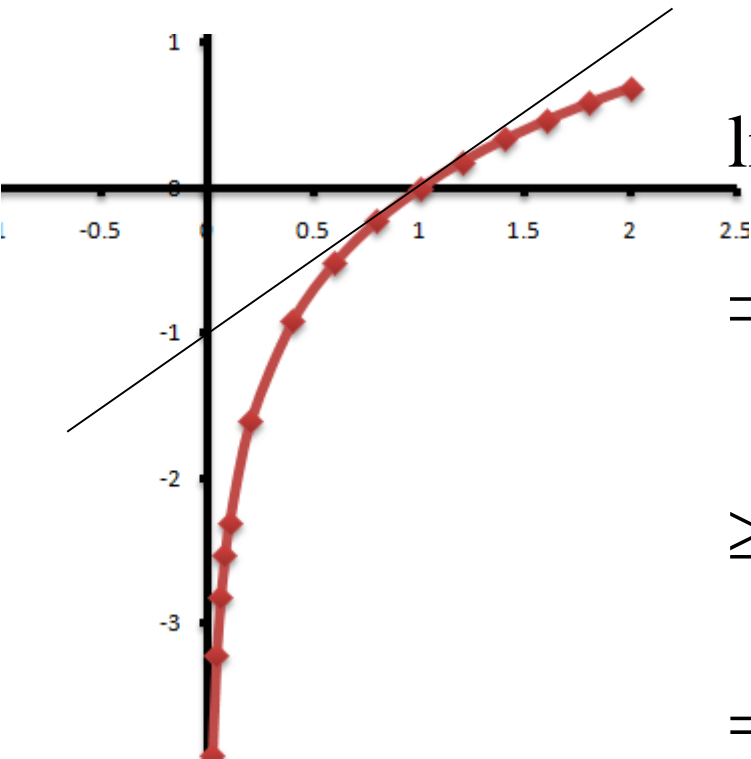
$$\ln x \leq x - 1 \Rightarrow$$

$$\Rightarrow \sum_{\pi} -P(\pi | s, \theta^0) \cdot \ln \frac{P(\pi | s, \theta)}{P(\pi | s, \theta^0)} \geq$$

$$\geq \sum_{\pi} -P(\pi | s, \theta^0) \cdot \left( \frac{P(\pi | s, \theta)}{P(\pi | s, \theta^0)} - 1 \right) =$$

$$= \sum_{\pi} [-P(\pi | s, \theta) + P(\pi | s, \theta^0)] =$$

$$= -\sum_{\pi} P(\pi | s, \theta) + \sum_{\pi} P(\pi | s, \theta^0) = 0$$



# WHY? Expectation-Maximisation algorithm

The EM algorithm is an iterative process

Each iteration performs two steps:

E-step: evaluation of  $Q(\theta | \theta^0) = \sum_{\pi} P(\pi | s, \theta^0) \cdot \log P(s, \pi | \theta)$

M-step: Maximisation of  $Q(\theta | \theta^0)$  over all  $\theta$

It does NOT assure to converge to the GLOBAL Maximum Likelihood

# Baum-Welch implementation of the EM algorithm

## *E-step:*

$$Q(\theta | \theta^0) = \sum_{\pi} P(\pi | s, \theta^0) \cdot \log P(s, \pi | \theta)$$

$$P(s, \pi | \theta) = a_{0, \pi(1)} \cdot \prod_{i=1}^T a_{\pi(i), \pi(i+1)} \cdot e_{\pi(i)}(s^i) =$$

$$= \prod_{k=0}^N \prod_{l=1}^N a_{k,l}^{A_{k,l}(\pi)} \cdot \prod_{k=1}^N \prod_{c \in \mathcal{C}} e_k(c)^{E_k(c, \pi)}$$

$A_{k,l}(\pi)$ : number of transitions between the states  $k$  and  $l$  in path  $\pi$

$E_k(c, \pi)$ : number of emissions of character  $c$  in path  $\pi$

$$A_{k,l} = \sum_{\pi} P(\pi | s, \theta^0) \cdot A_{k,l}(\pi)$$

Expected values over all the  
“actual” paths

$$E_k(c) = \sum_{\pi} P(\pi | s, \theta^0) \cdot E_k(c, \pi)$$

So:

$$Q(\theta | \theta^0) = \sum_{k=0}^N \sum_{l=1}^N A_{k,l} \cdot \ln a_{k,l} + \sum_{k=1}^N \sum_{c \in \mathcal{C}} E_k(c) \cdot \ln e_k(c)$$

# Baum-Welch implementation of the EM algorithm

*M-step:*

$$\frac{\partial Q}{\partial a_{k,l}} = 0 \quad \text{For any state } k \text{ and } l, \text{ with } \sum_l a_{k,l} = 1$$

$$\frac{\partial Q}{\partial e_k(c)} = 0 \quad \text{For any state } k \text{ and character } c, \text{ with } \sum_c e_k(c) = 1$$

By means of Lagrange's multipliers techniques, we can solve the system

$$L(a, e, \lambda, \mu) = \sum_{k=0}^N \sum_{l=1}^N A_{kl} \ln a_{kl} + \sum_{k=0}^N \sum_{c \in C} E_k(c) \ln e_k(c) - \sum_{i=1}^N \lambda_i \left( \sum_{j=1}^N a_{ij} - 1 \right) - \sum_{i=1}^N \mu_i \left( \sum_{c \in C} e_i(c) - 1 \right)$$

2 Lagrange's multipliers for each state are introduced

# Baum-Welch implementation of the EM algorithm

*M-step:*

$$L(a, e, \lambda, \mu) = \sum_{k=0}^N \sum_{l=1}^N A_{kl} \ln a_{kl} + \sum_{k=0}^N \sum_{c \in C} E_k(c) \ln e_k(c) - \sum_{i=1}^N \lambda_i \left( \sum_{j=1}^N a_{ij} - 1 \right) - \sum_{i=1}^N \mu_i \left( \sum_{c \in C} e_i(c) - 1 \right)$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial a_{kl}} = 0, \forall k, l \\ \frac{\partial L}{\partial e_k(c)} = 0, \forall k, c \\ \frac{\partial L}{\partial \lambda_k} = 0, \forall k \\ \frac{\partial L}{\partial \mu_k} = 0, \forall k \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \frac{A_{kl}}{a_{kl}} - \lambda_k = 0, \forall k, l \\ \frac{E_k(c)}{e_k(c)} - \mu_k = 0, \forall k, c \\ \sum_{j=1}^N a_{kj} = 1, \forall k \\ \sum_{c \in C} e_k(c) = 1, \forall k \end{array} \right.$$

# Baum-Welch implementation of the EM algorithm

*M-step:*

$$L(a, e, \lambda, \mu) = \sum_{k=0}^N \sum_{l=1}^N A_{kl} \ln a_{kl} + \sum_{k=0}^N \sum_{c \in C} E_k(c) \ln e_k(c) - \sum_{i=1}^N \lambda_i \left( \sum_{j=1}^N a_{ij} - 1 \right) - \sum_{i=1}^N \mu_i \left( \sum_{c \in C} e_i(c) - 1 \right)$$

$$\left\{ \begin{array}{l} \frac{A_{kl}}{a_{kl}} - \lambda_k = 0, \forall k, l \\ \frac{E_k(c)}{e_k(c)} - \mu_k = 0, \forall k, c \\ \sum_{j=1}^N a_{kj} = 1, \forall k \\ \sum_{c \in C} e_k(c) = 1, \forall k \end{array} \right. \Rightarrow \left\{ \begin{array}{l} a_{kl} = \frac{A_{kl}}{\lambda_k} \forall k, l \\ e_k(c) = \frac{E_k(c)}{\mu_k}, \forall k, c \\ \sum_{j=1}^N a_{kj} = 1, \forall k \\ \sum_{c \in C} e_k(c) = 1, \forall k \end{array} \right.$$

# Baum-Welch implementation of the EM algorithm

*M-step:*

$$a_{kl} = \frac{A_{kl}}{\sum_{i=1}^N A_{ki}}$$

$$e_k(c) = \frac{E_k(c)}{\sum_{d \in C} E_k(d)}$$

## Baum-Welch implementation of the EM algorithm

How to compute the expected number of transitions and emissions over all the paths

$$F_k(i) = P(s^1 s^2 s^3 \dots s^i, \pi(i) = k)$$

$$B_k(i) = P(s^{i+1} s^{i+2} s^{i+3} \dots s^T / \pi(i) = k)$$

$$A_{kl} = \sum_i P(\pi_i = k, \pi_{i+1} = l | s) = \frac{\sum_i F_k(i) \cdot a_{kl} \cdot e_l(s^{i+1}) \cdot B_l(i+1)}{P(s)}$$

$$E_k(c) = \sum_{s^i=c} P(s^i = c, \pi_i = k | s) = \frac{\sum_{s^i=c} F_k(i) \cdot B_k(i)}{P(s)}$$

This is for 1 training sequence: for general training with many sequences, just sum over them



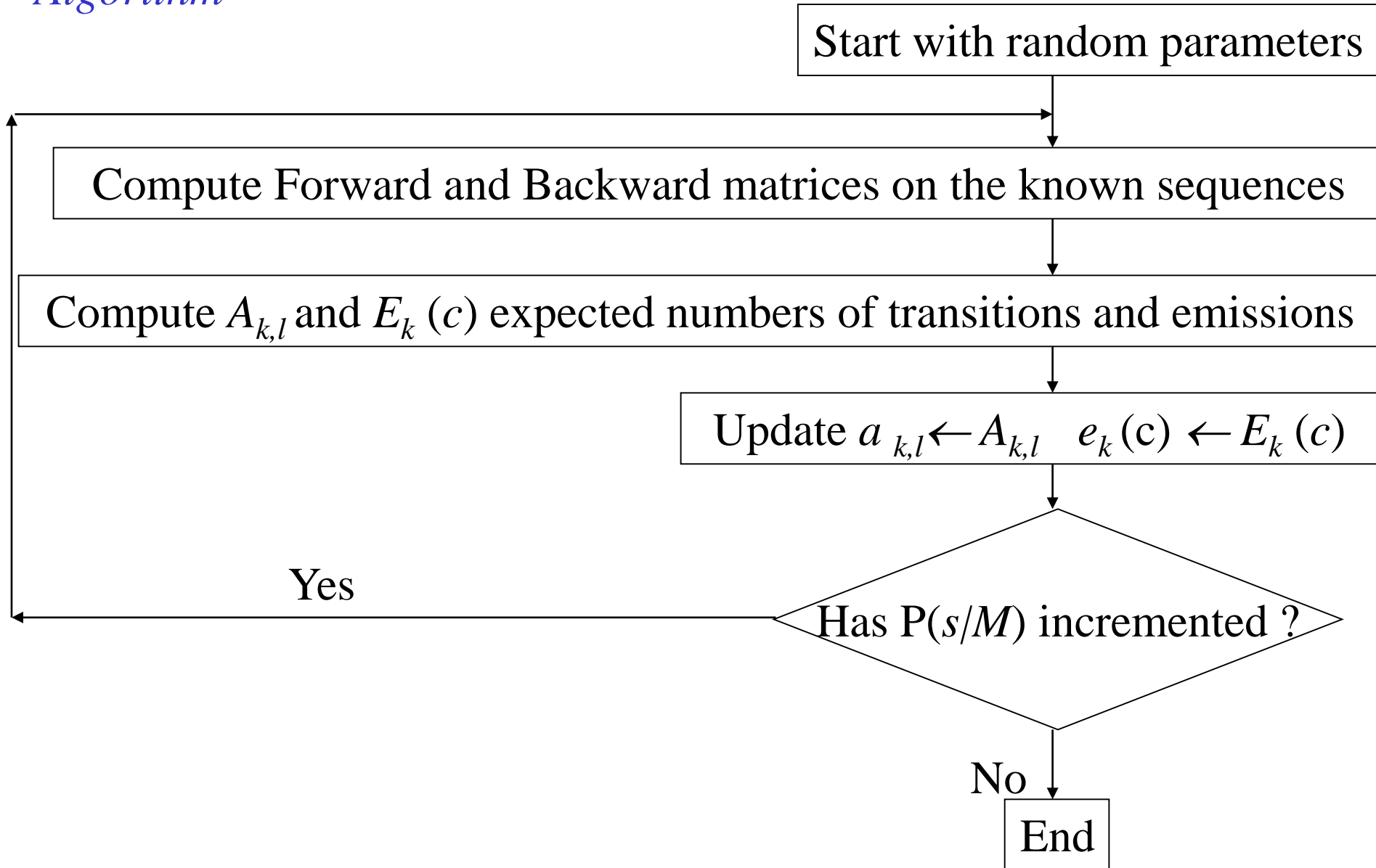
# Baum-Welch implementation of the EM algorithm: proof

$$\begin{aligned}
 A_{kl} &= \sum_i P(\pi_i = k, \pi_{i+1} = l \mid s) = \frac{P(s, \pi_i = k, \pi_{i+1} = l)}{P(s)} = \\
 &= \frac{\sum_i P(s^1, s^2 \dots s^i, \pi_i = k) \cdot P(\pi_{i+1} = l \mid \pi_i = k) \cdot P(s^{i+1} \mid \pi_{i+1} = l) \cdot P(s^{i+2}, \dots s^T \mid \pi_{i+1} = l)}{P(s)} = \\
 &= \frac{\sum_i F_k(i) \cdot a_{kl} \cdot e_l(s^{i+1}) \cdot B_l(i+1)}{P(s)}
 \end{aligned}$$

$$\begin{aligned}
 E_k(c) &= \sum_{s^i=c} P(s^i = c, \pi_i = k \mid s) = \frac{\sum_{s^i=c} P(s^1, s^2 \dots, s^i = c, s^{i+1}, \dots, s^T, \pi_i = k)}{P(s)} = \\
 &= \frac{\sum_{s^i=c} P(s^1, s^2 \dots, s^i = c, \pi_i = k) \cdot P(s^{i+1}, \dots, s^T \mid \pi_i = k)}{P(s)} = \frac{\sum_{s^i=c} F_k(i) \cdot B_k(i)}{P(s)}
 \end{aligned}$$

# Baum-Welch implementation of the EM algorithm

## *Algorithm*



# Profile HMMs

- HMMs for alignments

**PROLOGUE:**

**Pitfalls of standard alignments**

# Scoring a pairwise alignment

A: ALA**E**VLIRLIT**K**LYP  
 B: ASAK**K**HLNRLIT**E**LYP

$$Score(A, B) = \sum s(A^i, B^i)$$

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	-4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	-3
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	-3	-1	3
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1
K	-1	2	0	-1	-3	1	2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2	
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	2	-3
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	-1
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

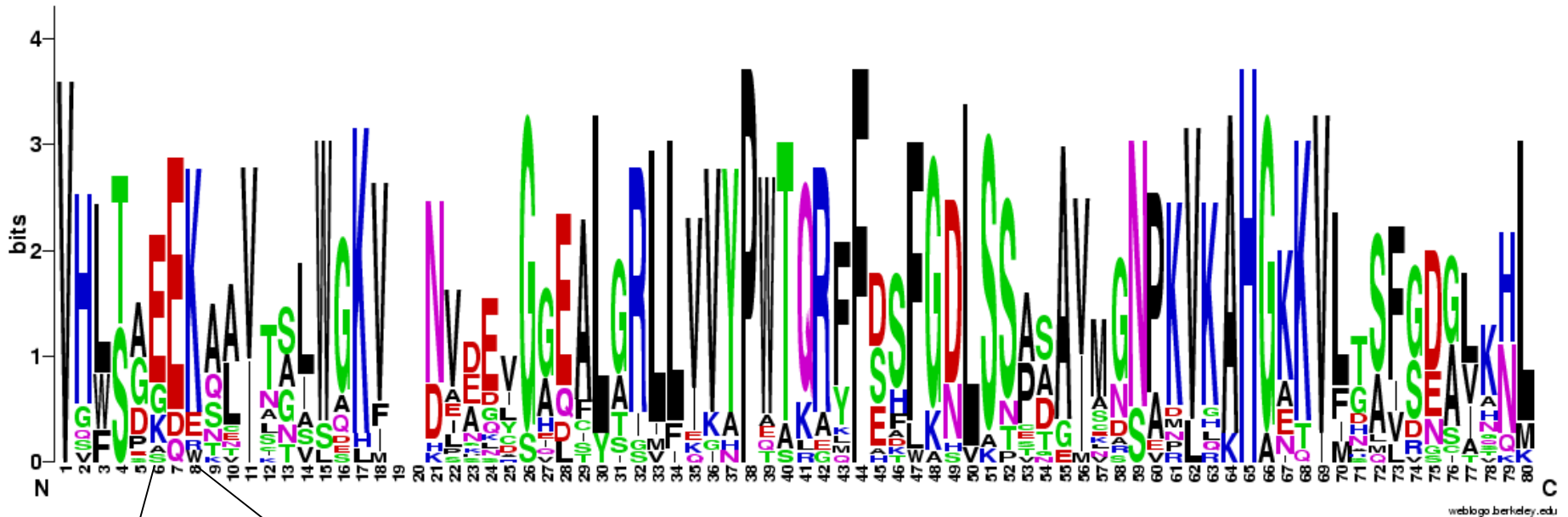
Blosum62

# Alignment of a family (globins)

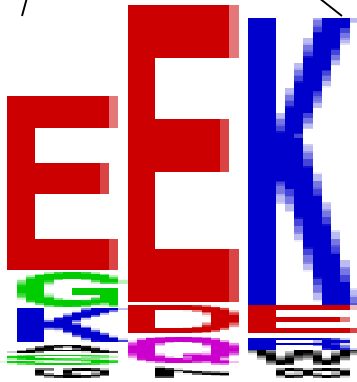
		10	20	30	40	50	60	70	80										
lqb1	pea/1-471	-GFTDKQ	EALVNSSSE	-FKQNL	PGYSILFY	TIVLEK	APAAKGL	FSFLKD	---	TAGVED	SPKLQ	AHAEQ	VFGLVR	DSSAAQ	L				
lqb1	vicfa/1-471	-GFTEKQ	EALVNSSSQ	LFKQNP	SNYSVL	FYTIIL	QKAPTAK	MFSLKD	---	SAGVVD	SPKLGA	HAKEK	VFGMVR	DSAV	QL				
hbb	speci/1-471	VHLS	DGEKNA	ISTAWG	KV--	HAAEV	GAEALGR	LLVVYP	PWTQRF	FDSFGD	LSASAV	MGNAKV	KAHGKK	VIDS	FSNGLKHL				
hbb	speto/1-471	VHLS	DGEKNA	ISTAWG	KV--	NAAEI	GAEALGR	LLVVYP	PWTQRF	FDSFGD	LSASAV	MGNAKV	KAHGKK	VIDS	FSNGLKHL				
hbb	equhe/1-471	VQLS	GEEKAA	AVLALW	DKV--	NEEEV	GGEALGR	LLVVYP	PWTQRF	FDSFGD	LSNPAA	VMGNP	KVKAH	GKKVLH	SFGE	GVHHL			
hbb	sunmu/1-471	VHLS	GEEKAC	VTGLWG	KV--	NEDEV	GAEALGR	LLVVYP	PWTQRF	FDSFGD	LSASAV	MGNPKV	KAHGKK	VLHSL	GEGVANL				
hbb	tupql/1-471	VHLS	GEEKAA	VTGLWG	KV--	DLEKV	GGQSLG	SLLIVY	PWTQRF	FDSFGD	LSSPSA	VMSNP	KVKAH	GKKVLT	SFSD	GDLNHL			
hbb	calar/1-471	VHLT	GEEKSA	VTALWG	KV--	NVDEV	GGEALGR	LLVVYP	PWTQRF	FESFGD	LSTPDA	VMNPN	KVKAH	GKKVLG	AFSD	GLTHL			
hbb	mansp/1-471	VHLT	PEEKTA	VTTLWG	KV--	NVDEV	GGEALGR	LLVVYP	PWTQRF	FDSFGD	LSSPDA	VMGNP	KVKAH	GKKVLG	AFSD	GDLNHL			
hbb	rabit/1-471	VHLS	SSEKSA	VTALWG	KV--	NVEEV	GGEALGR	LLVVYP	PWTQRF	FESFGD	LSANAV	MNPNK	VKAH	GKKVLA	AFSE	GLSHL			
hbb	ursma/1-471	VHLT	GEEKSL	VTLWG	KV--	NVDEV	GGEALGR	LLVVYP	PWTQRF	FDSFGD	LSADAI	MNPNK	VKAH	GKKVLN	SFSD	GDLKNL			
hbb	triin/1-471	VHLT	PEEKAL	VIGLWA	KV--	NVKEY	GGEALGR	LLVVYP	PWTQRF	FEHFGD	LSASAI	MNPNK	VKAH	GEKVFT	SFGD	GLKHL			
hbb	ornan/1-471	VHLS	GGEKSA	VTNLWG	KV--	NINEL	GGEALGR	LLVVYP	PWTQRF	FEAFGD	LSAGAV	MGNPK	VKAH	GAKVLT	SFGD	ALKNL			
hbb	tacac/1-471	VHLS	GSEKTA	VTNLWG	HV--	NVNEL	GGEALGR	LLVVYP	PWTQRF	FESFGD	LSADAV	MGNAK	VKAH	GAKVLT	SFGD	ALKNL			
hbe	ponpy/1-471	VHFT	AEEKAA	VTSLWS	KM--	NVEE	AGGEALGR	LLVVYP	PWTQRF	FDSFGN	LSSPSA	ILGNP	KVKAH	GKKVLT	SFGD	AIKNM			
hbb	colli/1-471	VHWS	AEEKQL	ITSIWG	KV--	NVAD	CGAEAL	ARLLIV	YPWTQ	RFFSS	FGNLSS	ATAIS	GNPNV	KAH	GKKVLT	SFGD	AVKNL		
hbb	larri/1-471	VHWS	AEEKQL	ITGLWG	KV--	NVAD	CGAEAL	ARLLIV	YPWTQ	RFFAS	FGNLSS	PTAIN	GNPMV	RAH	GKKVLT	SFGE	AVKNL		
hbb1	varex/1-471	VHWT	AEEKQL	ICSLWG	KI--	DVGL	IGGETL	AGLLVI	YPWTQ	RQFSH	FGNLSS	PTAIA	GNPRV	KAH	GKKVLT	SFGD	AIKNL		
hbb2	xentr/1-471	VHWT	AEEKAT	IASVWG	KV--	DIEQ	DGHDAL	SRLLV	VYPWTQ	RYFSS	FGNLSS	NVSAV	SGNVK	KAH	GNKVL	SAVG	SAIQHL		
hbb1	ranca/1-471	VHWT	AEEKAV	INSVWQ	KV--	DVEQ	DGHEAL	TRLFIV	YPWTQ	RYFST	FGDLSS	PAIAI	GNPKV	HAH	GKKIL	GAID	NAIHNL		
hbb2	tracr/1-471	VHLT	AE	DRKEIA	AAILG	KV--	NVDS	LGGOCL	ARLI	VNPWS	RRYFH	DFGDL	SSCDA	ICRN	PKVL	AHGAK	VMRS	I	VEATKHL
hba4	salir/1-471	-SL	SAKDKAN	VKAIW	GKILPK	SDEIGE	QALS	RMLV	VYPQTK	AYFSH	MASVAP	----	GSAP	VKKH	GITIM	NQID	DCVGHM	:	
myg	escqi/1-471	-VL	SDAEW	QLVLN	IAKVE	ADVAG	HGQD	ILIR	LFKG	HPETL	EKFDK	FKHLK	TEAEM	KASE	DLKKH	GNTVL	TALG	GILKKK	:

Different positions are not equivalent

# Sequence logos



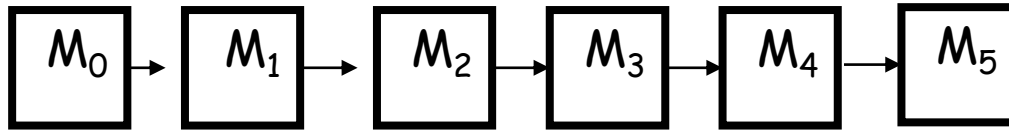
<http://weblogo.berkeley.edu/cache/file5h2DWc.png>



The substitution score IN A FAMILY should depend on the position (the same for gaps)

For modelling families we need more flexible tools

# How to align?



Each state represent a position in the alignment.

A	C	G	G	T	A
$M_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$

A	C	G	A	T	C
$M_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$

A	T	G	T	T	C
$M_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$

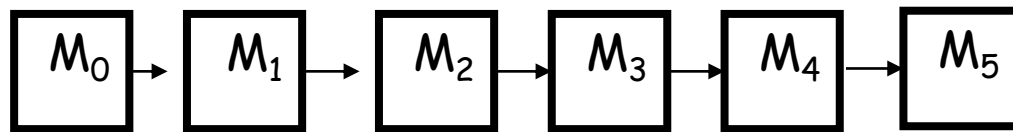
Each position has a peculiar composition



Given a set of sequences..

A	C	G	G	T	A
A	C	G	A	T	C
A	T	G	T	T	C

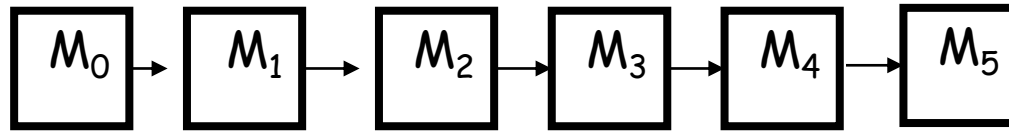
..we can train a model..



..estimating the emission probabilities.

A	1	0	0	0.33	0	0.33
C	0	0.66	0	0	0	0.66
G	0	0	1	0.33	0	0
T	0	0.33	0	0.33	1	0

Given a trained model..



A	1	0	0	0.33	0	0.33
C	0	0.66	0	0	0	0.66
G	0	0	1	0.33	0	0
T	0	0.33	0	0.33	1	0

..we can align a new sequence..

A      C      G      A      T      C

..computing the probability of generating it

$$P(s|M) = 1 \times 0.66 \times 1 \times 0.33 \times 1 \times 0.66$$

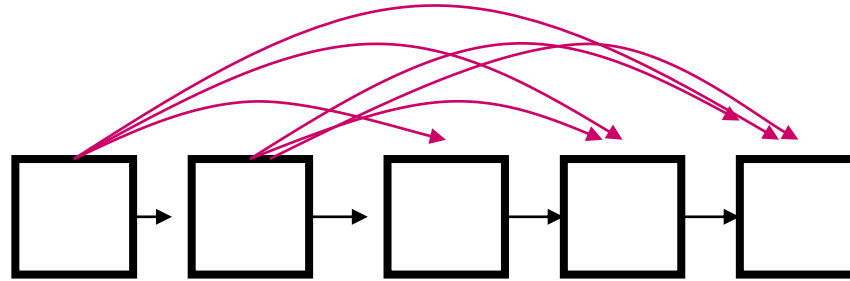
And for the sequence AGATC ?

	$M_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$
A	1	0	0	0.33	0	0.33
C	0	0.66	0	0	0	0.66
G	0	0	1	0.33	0	0
T	0	0.33	0	0.33	1	0

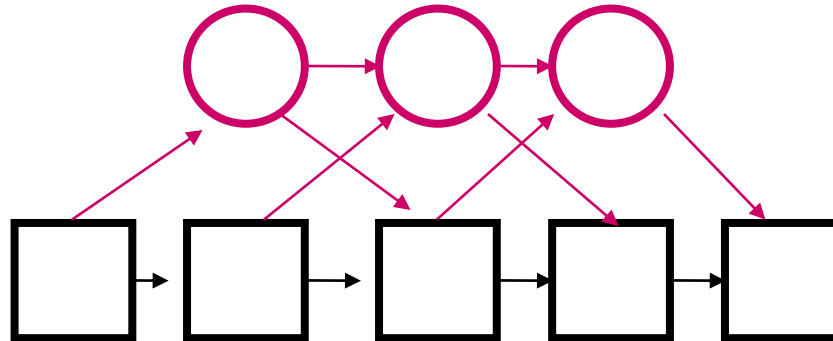
A		G	A	T	C
$M_0$	$M_2$	$M_3$	$M_4$	$M_5$	$M_5$

We need a way to introduce gaps

## Silent states

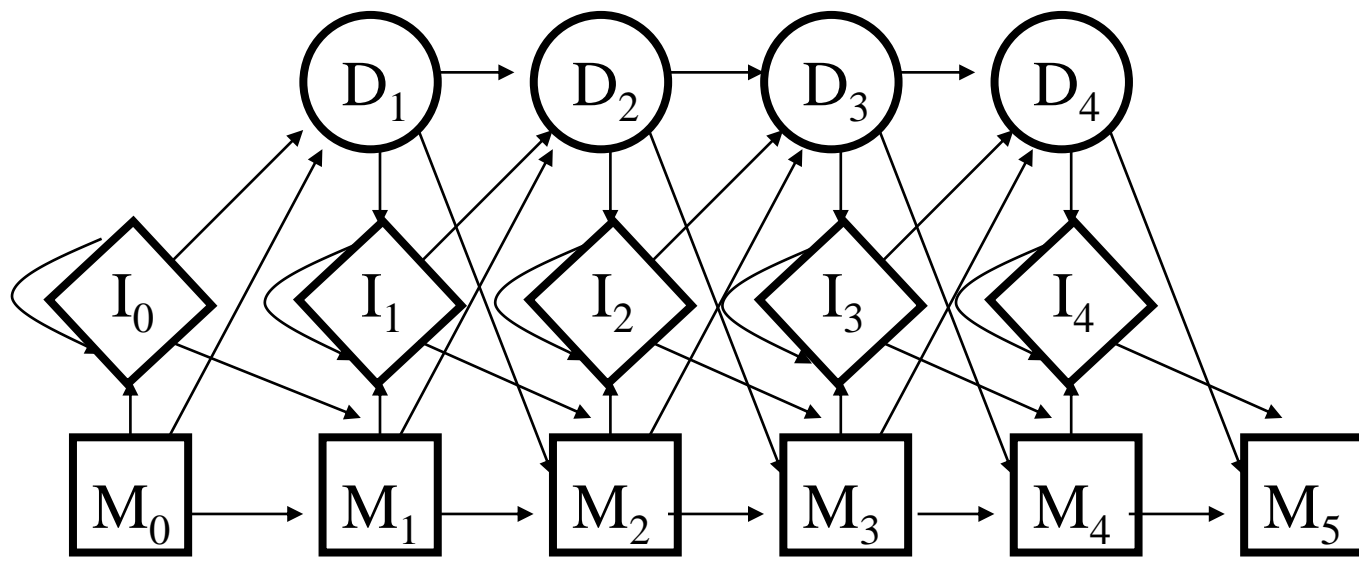


Red transitions allow gaps  
 $N(N-1)/2$  transitions



To reduce the number of parameters we can use states that  
**doesn't emit any character**  
 $4N-8$  transitions

# Profile HMMs



Delete states

Insert states

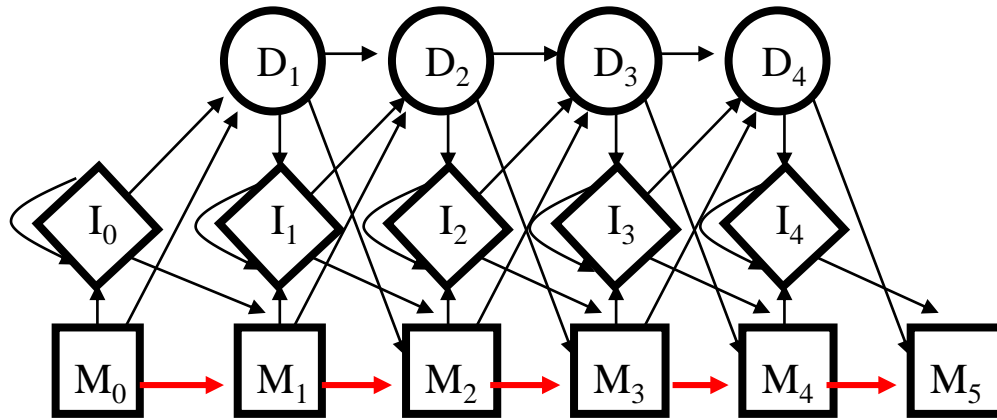
Match states

<b>A</b>	<b>C</b>	<b>G</b>	<b>G</b>	<b>T</b>	<b>A</b>
$M_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$

<b>A</b>	<b>C</b>	<b>G</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>C</b>
$M_0$	$I_0$	$I_0$	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$

<b>A</b>	<b>G</b>	<b>A</b>	<b>T</b>	<b>C</b>
$M_0$	$D_1$	$M_2$	$M_3$	$M_5$

# Example of alignment



*Sequence 1*

**A S T R A L**

*Viterbi path*

**$M_0$   $M_1$   $M_2$   $M_3$   $M_4$   $M_5$**

**A S T R A L**

*Sequence 2*

**A S T A I L**

*Viterbi path*

**$M_0$   $M_1$   $M_2$   $D_3$   $M_4$   $I_4$   $M_5$**

**A S T A I L**

*Sequence 3*

**A R T I**

*Viterbi path*

**$M_0$   $M_1$   $M_2$   $D_3$   $D_4$   $M_5$**

**A R T I**

# Example of alignment

$M_0$	$M_1$	$M_2$	$M_3$	$M_4$		$M_5$	
A	S	T	R	A		L	<i>Sequence 1</i>
$M_0$	$M_1$	$M_2$	$D_3$	$M_4$	$I_4$	$M_5$	
A	S	T		A	I	L	<i>Sequence 2</i>
$M_0$	$M_1$	$M_2$	$D_3$	$D_4$		$M_5$	
A	R	T				I	<i>Sequence 3</i>

Grouping by vertical layers

	0	1	2	3	4	5
$s_1$	A	S	T	R	A	L
$s_2$	A	S	T		AI	L
$s_3$	A	R	T			I

**Alignment**

ASTRA-L

AST-AIL

ART---I

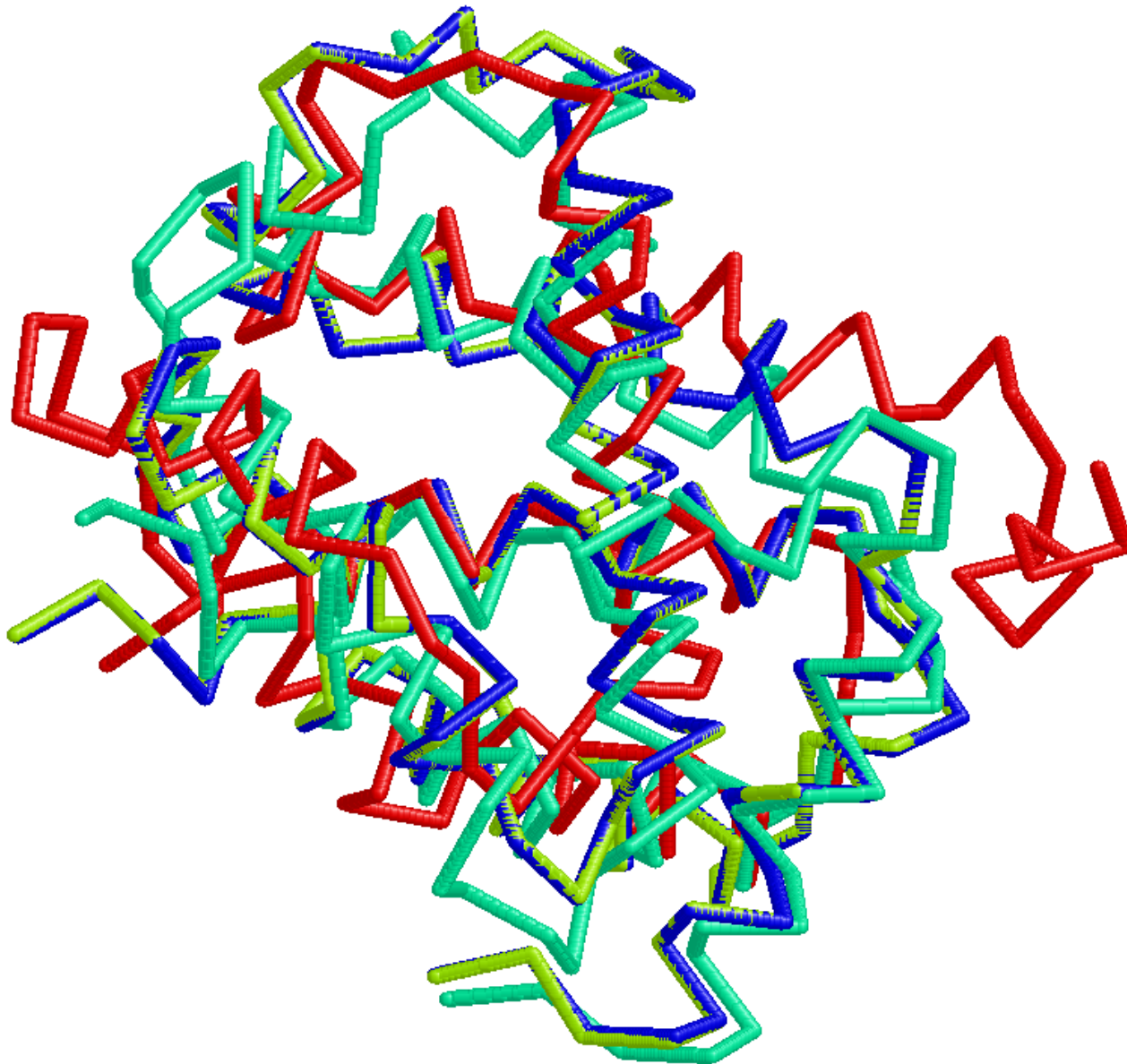
$-\text{Log } P(s \mid M)$  Is an alignment score

# Profile HMMs

- HMMs for alignments
- Example on globins



# Structural alignment of globins



# Structural alignment of globins

```
AAAAAAAAAAAAAAAAA   BBBBCCCCCCCCCCCC
                                DDDD
-----VLSPADKTNVKAANGKVGA--HAGEYGAEALERMFLSFPTTKTYFPHF-DL
-----VHLTPEEKSAVTALWGKV---NVDEVGGEALGRLLVVYPWTQRFFESFGDL
-----VLSEGEWQLVLHVWAKVEA--DIAGHGQDILIRLFKHHPETLEKFDREFKHL
-----LSADQISTVQASFDKVKG-----DPVGILYAVFKADPSIMAKFTQFAG-
PIVDTGSAVPLSAAEKTKIRSAWAPVYS--TYETSGVDILVKFFTSTPAAQEFFPKFKGL
-----GALTESQAALVKSSWEEFN--NIPKHTRFFILVLEIAPAAKDLFS-FLK-
-----GLSAAQRQVIAATWKDIAGADNGAGVGKDCLIKFLSAHPQMAAVFG-FSG-
```

```
DDDDDDDEEEEEEEEEEEEEEEEEEEEEEEEEEE          FFFFFFFF   FFGGG
                                F           GG   GG
S-----HGSAQVKGHGKKVADALTNAVAHV--D--DMPNALSALSDLHAHKL--RVDPV
STPDVAVMGNPKVKAHGKKVLGAFSDGLAHL---D--NLKGTfATLSELHCDKL--HVDPE
KSEAEMKASEDLKKHGVTVLTAAILKK---K-GHHEAELKPLAQSHATKH--KIPIK
KDLESIKGTAPFETHANRIVGFFSKIIGEL--P---NIEADVNTFVASHKPRG---VTHD
TTADQLKKSADVRWHAERIINAVNDAVASM--DTEKMSMKLRDLGSKHAKSF--QVDPQ
GTSEVPQNNPELQAHAGKVFKLVYEAIIQLQVTGVVVTDATLKNLGSVHVSKG---VADA
---AS---DPGVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKAVGVRHKG YGNKHIKAQ
```

```
GGGGGGGGGGGGGGGGGG   HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
NFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR-----
NFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH-----
YLEFISEAIIHVLHSRHPADFGADAQGAMSKALELFRKDIAAKYKELGYQG
QLNNFRAGFVSYMKAHT--DFA-GAEAAWGATLDTFFGMIFSKM-----
YFKVLA AVIADTVAAG-----DAGFEKLMSMICILLRSAY-----
HFPVVK EAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA--
YFEPLGASLLSAMEHRIGGKMNAAKDAWAAAYADISGALISGLQS-----
```

Bashdorf D, Chothia C & Lesk AM, (1987) *Determinants of a protein fold: unique features of the globin amino sequence*. **J.Mol.Biol.** 196, 199-216

# Alignment of globins reconstructed with profile HMMs

The parameters of an HMM are estimated from a training set of 400 unaligned sequences. After the HMM is built, it is used to obtain a multiple alignment of all the training sequences. This is the alignment of the 7 globins as aligned with the trained model

```

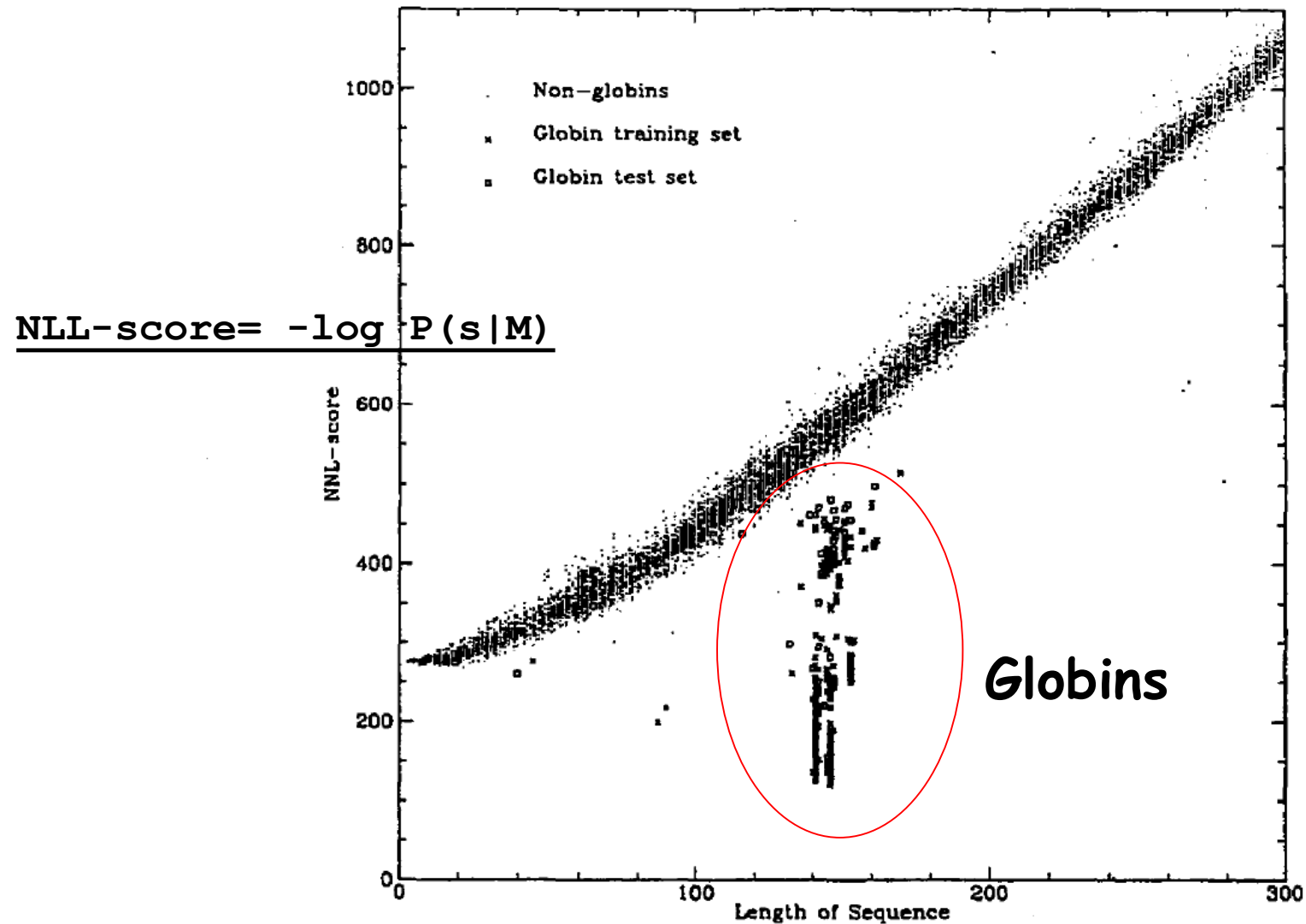
      AAAAAAAAAAAAAAAAAA  BBBB BBBB BBBB BBBB BBBBBB CCCCCCCCCCCC
                                DDDD
      *****
V.....LSPADKTNVKA AWGKVGA..HAGEYGAEALERMF LSFPTTKTYFPHFD-L
Vh.....LTPEEKSAVTALWGKV--.NVDEVGGEALGRLLVVYPWTQRFFESFGDL
V.....LSEGEWQLVLHVWAKVEA..DVAGHGQDILIRLFKSHPETLEKFD RFKHL
-.....LSADQISTVQASFDKV--.KGDPVGI--LYAVFKADPSIMAKFTQFAGK
PivdtgsvapLSAAEKT KIRSAWAPVYS..TYETSGVDILVKFFTSTPAAQEFFPKFKGL
Ga.....LTESQAALVKSSWEEFNA..NIPKHTHRFFILVLEIAPAAKDLFSFLK-G
G.....LSAAQRQVIAATWKDIAGadNGAGVGKDCLIKFLSAHPQMA--AVFG-F

      DDDDDDDDEE  EEEEEEEEEEEEEEEEEEEEEEE
                                F
      *****
SHGSAQVKGH-GKK.---VADAL TNAVAHVDD....MPNALSALSDLHA...HKLRVD
STPDAVMGNPKVKA.HGKKVLGAFSDGLAHLDN....LKGTFATLSELHC...DKLHVD
KTEA-EMKASEDLKkHGVTVLTA LGAILKKKGH....HEAELKPLAQSHA...TKHKIP
DLES-IKGTAPFET.HANRIVGFFSKIIGELPN....IEADVNTFVASHK...PR-GVT
TTADQLKKSADVRW.HAERIINAVNDAVASMDDtek..MSMKLRDLSGKHA...KSFQVD
TSEVPQ-NNPELQA.HAGKVFKL VYEAAIQLQVtgvvvTDATLKNLGSVHV...SK-GVA
SGAS---DPGVAA.LGAKVLAQIGVAVSHLGDegk..MVAQMKA VGRHKGygNK-HIK

      GGGGGGGGGGGGGGGGGGGG  HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
      *****
PVNFKLLSHCLLVTLAAHLP AEFTPAVHASLDKFLASVSTVLTSKY.....R
PENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKV VAGVANALAHKY.....H
IKYLEFISEAIIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYkelgyqG
HDQLNNFRAGFVSYMKAH--TDF-AGAEAAWGATLD TFFGMIFSKM.....-
PQYFKVLA AVIADTVAA--GD-----AGFEKLMSMICILLRSAY.....-
DAHFPVVKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMnda...A
AQYFEPLGASLLSAMEHRIGGKMNA AAKDAWAAAYADISGALISGLq.....S
```

[Krogh A, Brown M, Mian IS, Sjolander K & Haussler D \(1994\) \*Hidden Markov Models in computational biology: applications to protein modelling\*. \*\*J.Mol.Biol.\*\* \*\*235\*\*, 1501-1531](#)

# Discrimination power of profile HMMs



**Figure 6.** Plot of NLL-score *versus* sequence length for globins and non-globins. All sequences of length less than 300 from the SWISS-PROT 22 database are shown, including partial sequences and 3 false globins from the globin file, and sequences from the database containing many Xs.

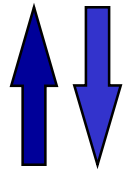
[Krogh A, Brown M, Mian IS, Sjolander K & Haussler D \(1994\) \*Hidden Markov Models in computational biology: applications to protein modelling\*. \*J.Mol.Biol.\* 235, 1501-1531](#)

# HMMs for Mapping problems

- Mapping problems in protein prediction

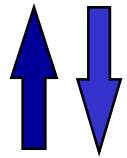
# Secondary structure

## Covalent structure

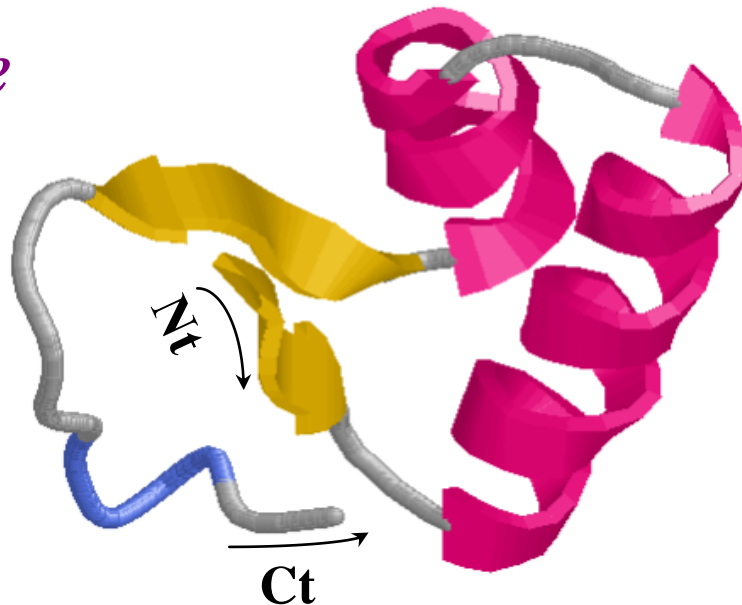


TTCCPSIVARSNFNVCRLPGTPEAICATYTGCIIPGATCPGDYAN

## Secondary structure



## 3D structure



# Topology of membrane proteins

## Topography

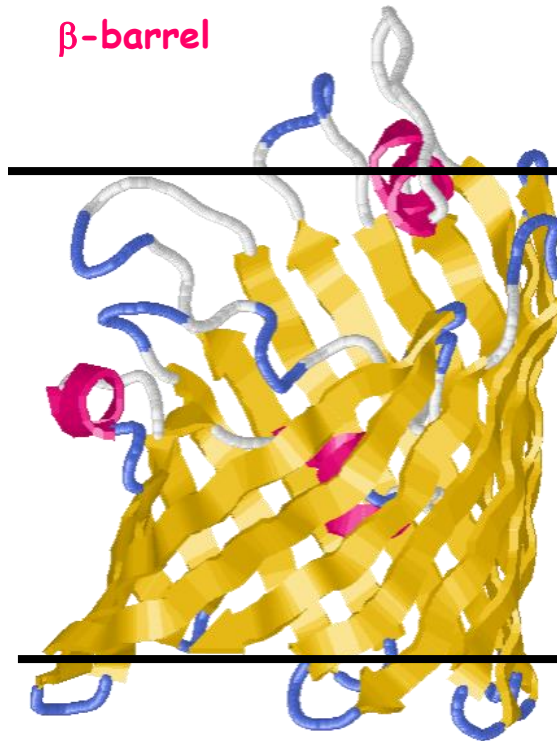
*position of Trans Membrane Segments along the sequence*

ALALMLCMLTYRHKELKCLKLKK ALALMLCMLTYRHKELKCLKLKK ALALMLCMLTYRHKELKCLKLKK



Outer Membrane

$\beta$ -barrel

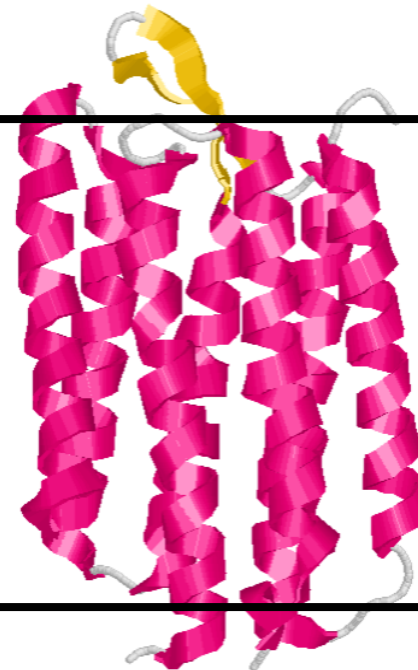


Porin  
(*Rhodobacter capsulatus*)

Inner Membrane

$\alpha$ -helices

Bilayer



Bacteriorhodopsin  
(*Halobacterium salinarum*)

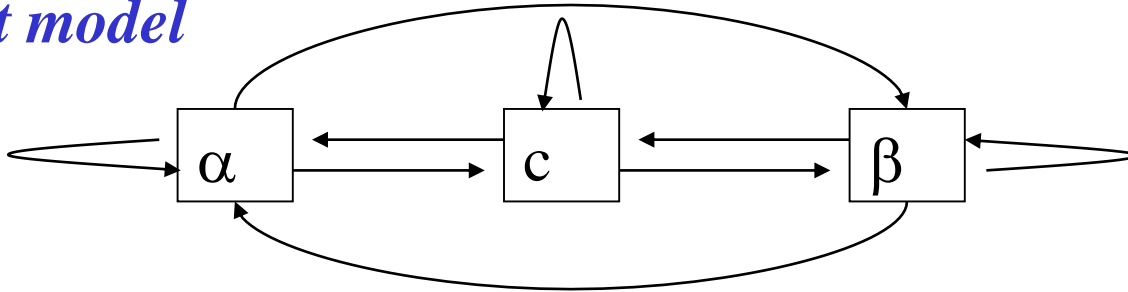
# HMMs for Mapping problems

- Mapping problems in protein prediction
- Labelled HMMs

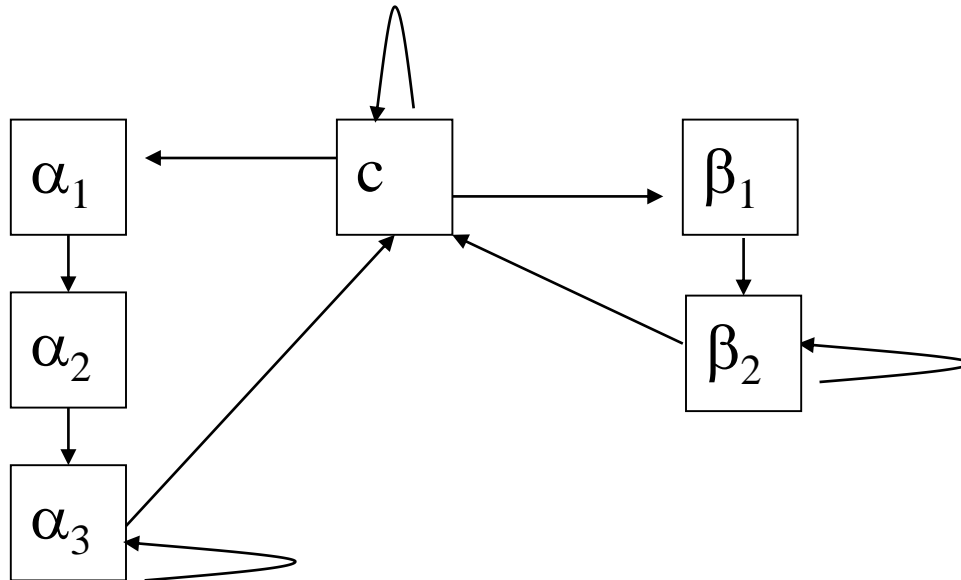


# HMM for secondary structure prediction

## *Simplest model*



## *Introducing a grammar*



# HMM for secondary structure prediction

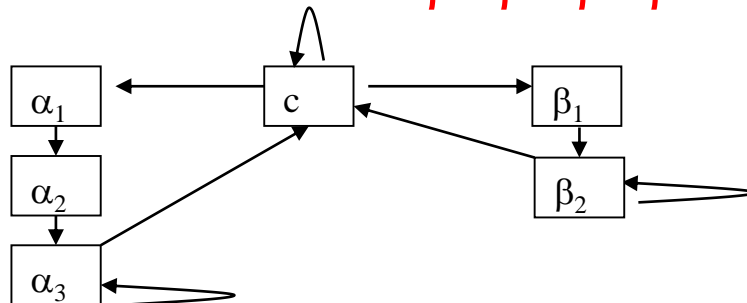
## Labels

The states  $\alpha_1$ ,  $\alpha_2$  and  $\alpha_3$  share the same *label*, so states  $\beta_1$  and  $\beta_2$  do.

Decoding the

Viterbi path for emitting a sequence  $s$ , makes a mapping between the sequence  $s$  and a sequence of labels  $y$

<b>S</b>	<b>A</b>	<b>L</b>	<b>K</b>	<b>M</b>	<b>N</b>	<b>Y</b>	<b>T</b>	<b>R</b>	<b>E</b>	<b>I</b>	<b>M</b>	<b>V</b>	<b>A</b>	<b>S</b>	<b>N</b>	<b>Q</b>	$s$ : <i>Sequence</i>
c	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_4$	$\alpha_4$	c	c	c	c	$\beta_1$	$\beta_2$	$\beta_2$	$\beta_2$	c	c	$\pi$ : <i>Path</i>
c	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$	c	c	c	c	$\beta$	$\beta$	$\beta$	$\beta$	c	c	$Y(\pi)$ : <i>Labels</i>



## Computing $P(s, y \mid M)$

$$P(s, y \mid M) = \sum_{\pi | Y(\pi) = y} P(s, \pi \mid M)$$

Only the path whose labelling is  $y$  have to be considered in the sum  
In Forward and Backward algorithms it means to set

$$F_k(i) = 0, \quad B_k(i) = 0 \quad \text{if} \quad Y(k) \neq y^i$$

		S	A	L	K	M	N	Y	T	R	E	I	M	V	A	S	N	Q	<i>s: Sequence</i> <i>y: Labels</i>
		c	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$	c	c	c	c	$\beta$	$\beta$	$\beta$	$\beta$	c	c	
$\alpha_1$	$\alpha$																		
$\alpha_2$	$\alpha$																		
$\alpha_3$	$\alpha$																		
$\alpha_4$	$\alpha$																		
$\beta_1$	$\beta$																		
$\beta_2$	$\beta$																		
c	c																		

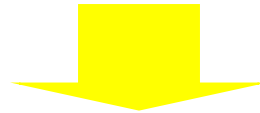
*States Labelling*

## Baum-Welch training algorithm for labelled HMMs

Given a set of known labelled sequences (e.g. amino acid sequences and their native secondary structure) we want to find the parameters of the model, without knowing the generating paths:

$$\theta^{\text{ML}} = \operatorname{argmax}_{\theta} [\text{P} ( s, y \mid \theta, M )]$$

The algorithm is the same as in the non-labelled case if we use the Forward and Backward matrices defined in the last slide.

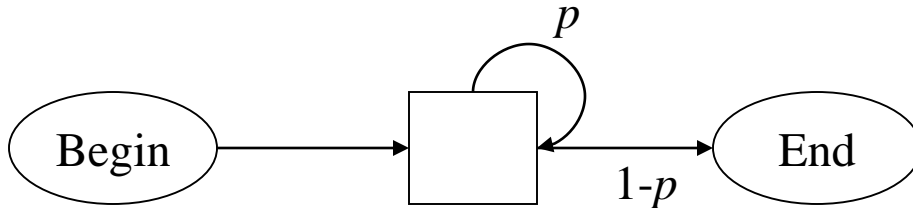


Supervised learning of the mapping

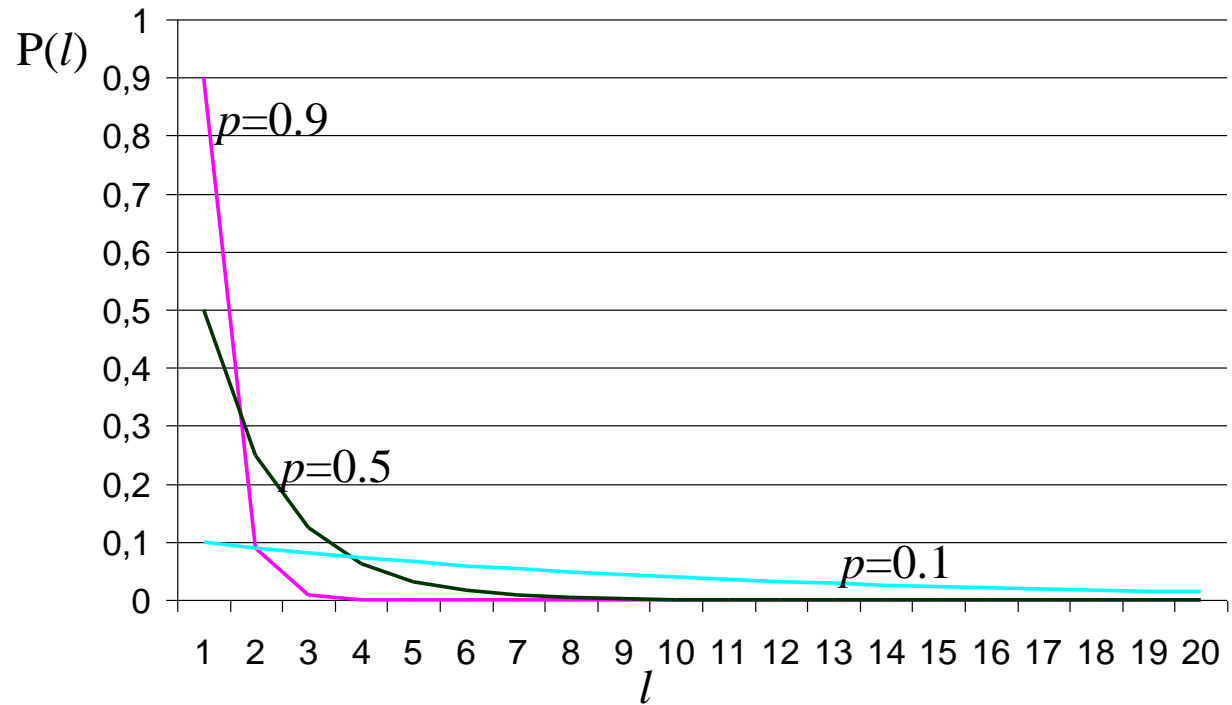
# HMMs for Mapping problems

- Mapping problems in protein prediction
- Labelled HMMs
- Duration modelling

# Self loops and geometric decay



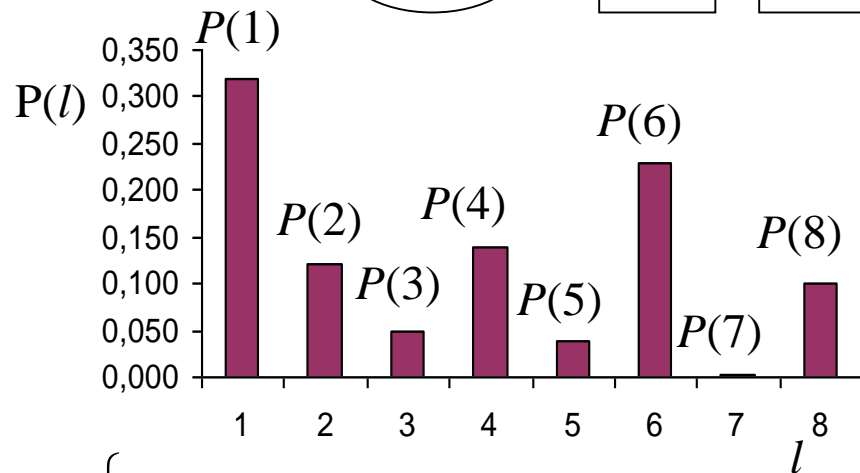
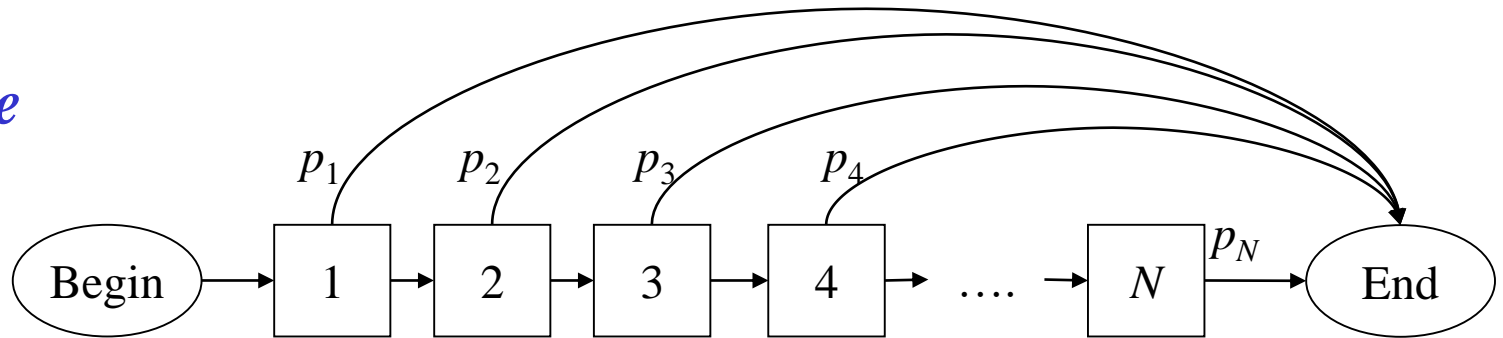
$$P(l) = p^{l-1} \cdot (1-p)$$



The length distribution of the generated segments is always exp-like

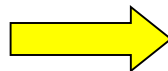
# How can we model other length distributions?

## Limited case



This topology can model any length distribution between 1 and  $N$

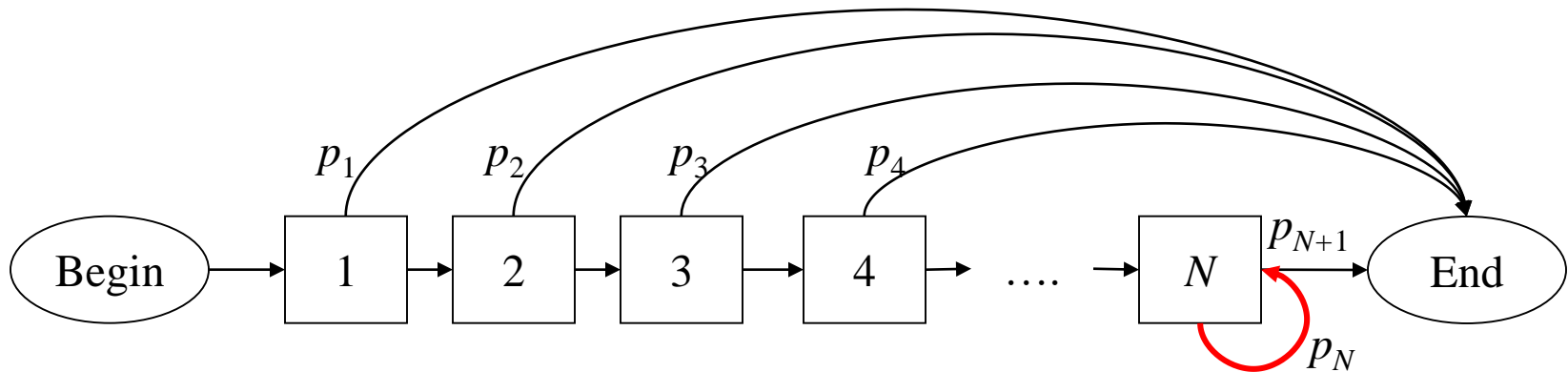
$$\left\{ \begin{array}{l} P(1) = p_1 \\ P(2) = (1 - p_1) \cdot p_2 \\ P(3) = (1 - p_1) \cdot (1 - p_2) \cdot p_3 \\ \dots\dots\dots \\ P(N) = \sum_{i=1}^{N-1} (1 - p_i) \cdot p_N \end{array} \right.$$



$$\left\{ \begin{array}{l} p_1 = P(1) \\ p_2 = P(2) / (1 - p_1) \\ \dots\dots\dots \\ p_k = \frac{P(k)}{\sum_{i=1}^{k-1} (1 - p_i)} \end{array} \right.$$

# How can we model other length distributions?

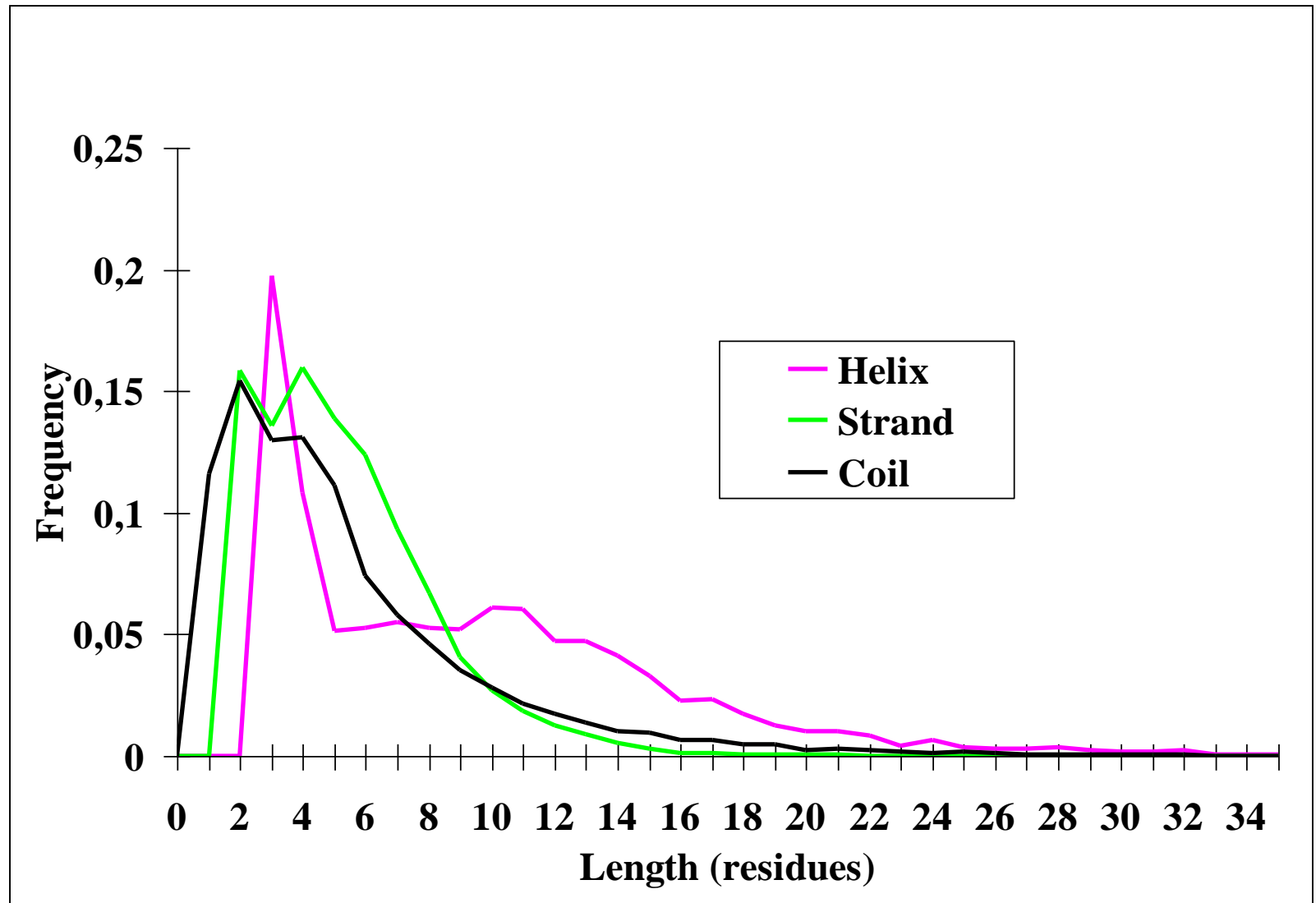
## *Non limited case*



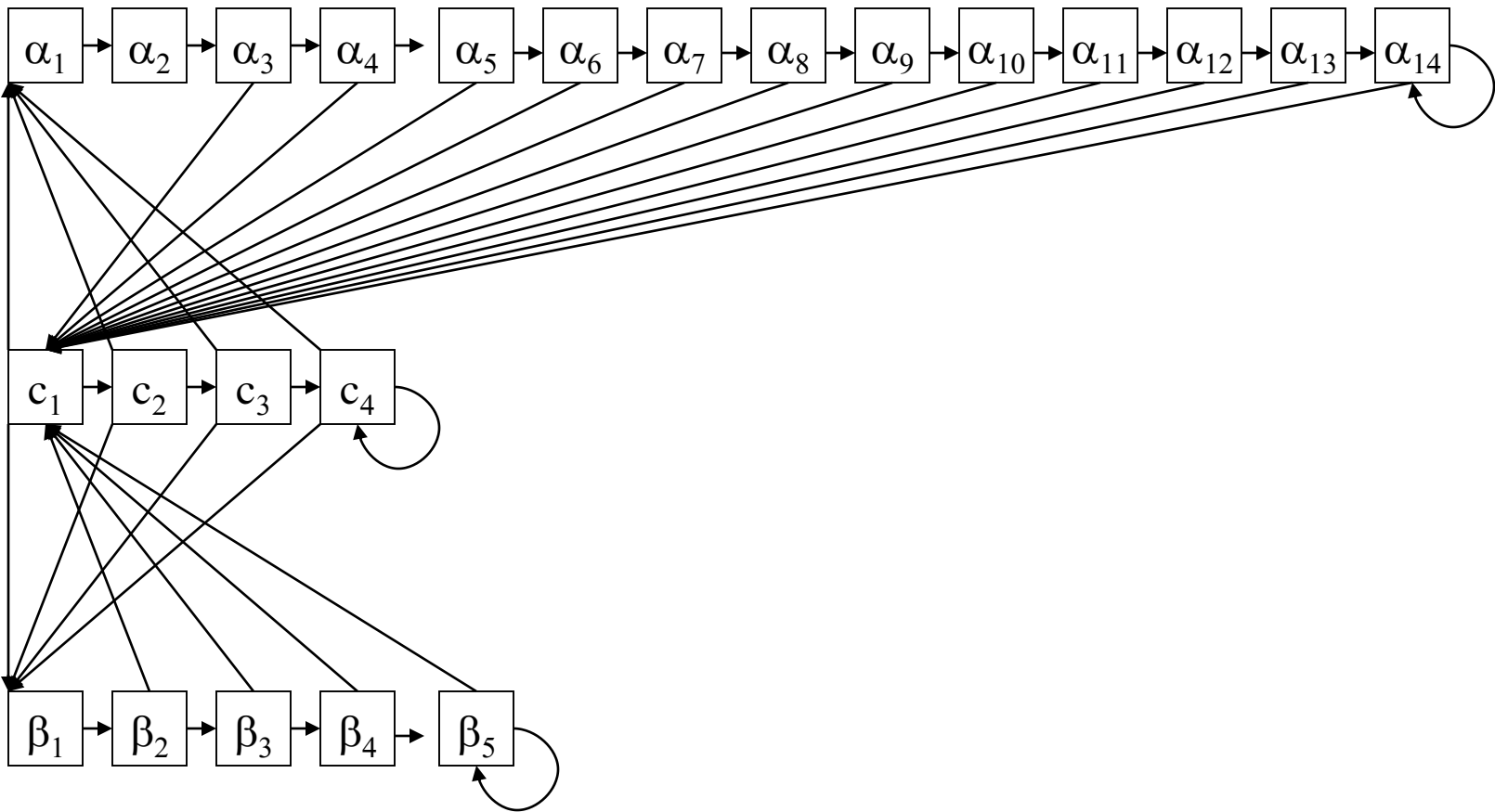
This topology can model any length distribution between 1 and  $N-1$  and a geometrical decay from  $N$  and  $\infty$



## Secondary structure: length statistic



## Secondary structure: model



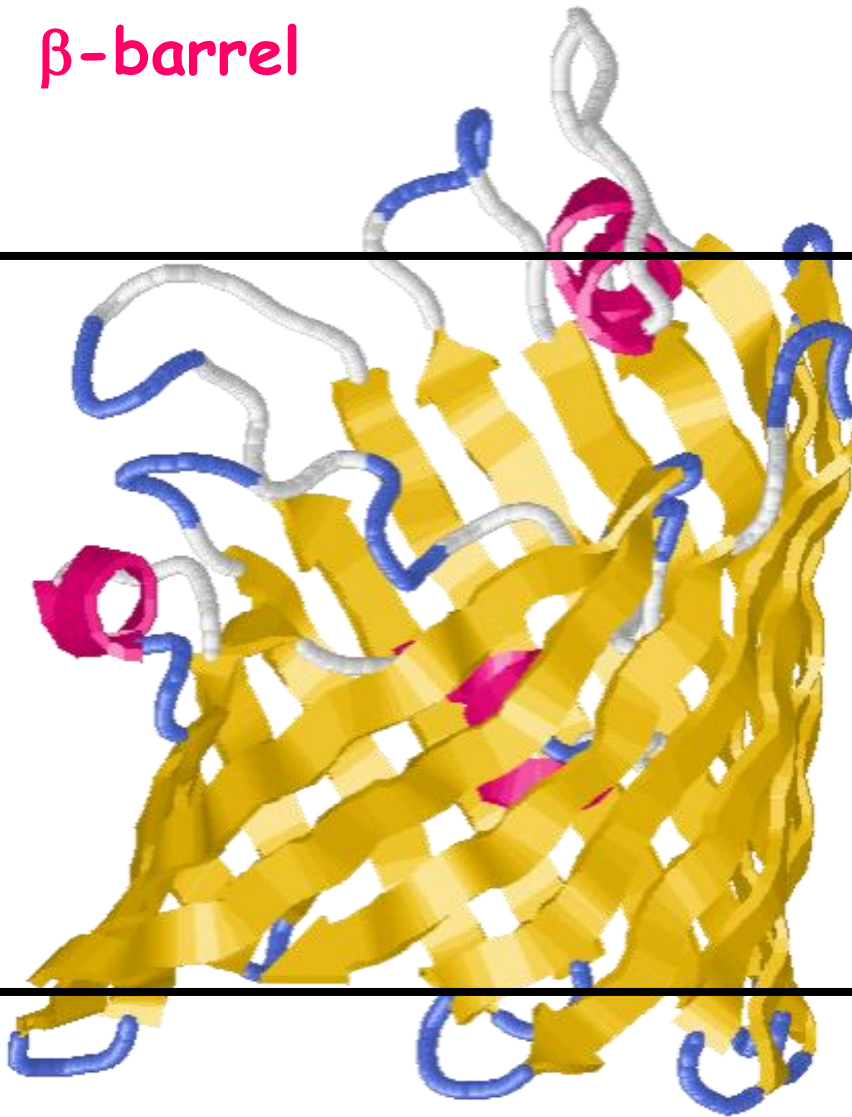
Do we use the same emission probabilities for states sharing the same label?

# HMMs for Mapping problems

- Mapping problems in protein prediction
- Labelled HMMs
- Duration modelling
- Models for membrane proteins

## Outer Membrane

$\beta$ -barrel



**Porin**

(*Rhodobacter capsulatus*)

## Inner Membrane

$\alpha$ -helices



**Bacteriorhodopsin**

(*Halobacterium salinarum*)

**Bilayer**

# Topology of membrane proteins

## Topography

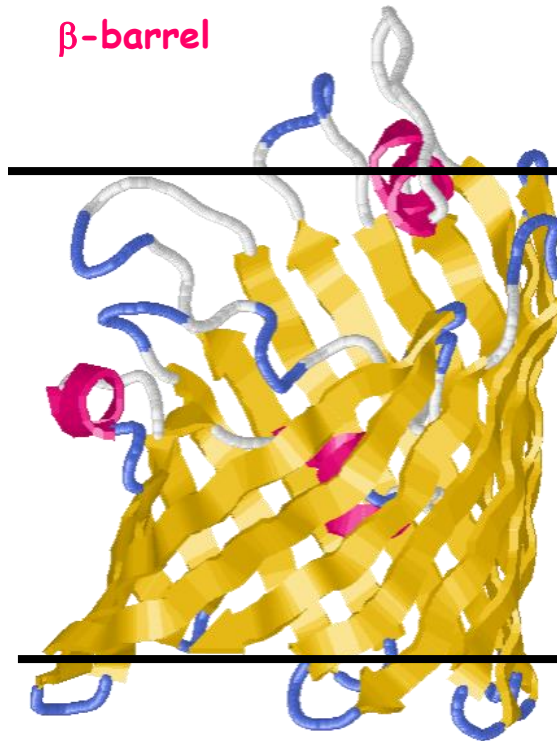
*position of Trans Membrane Segments along the sequence*

ALALMLCMLTYRHKELKCLKLKK ALALMLCMLTYRHKELKCLKLKK ALALMLCMLTYRHKELKCLKLKK



Outer Membrane

$\beta$ -barrel

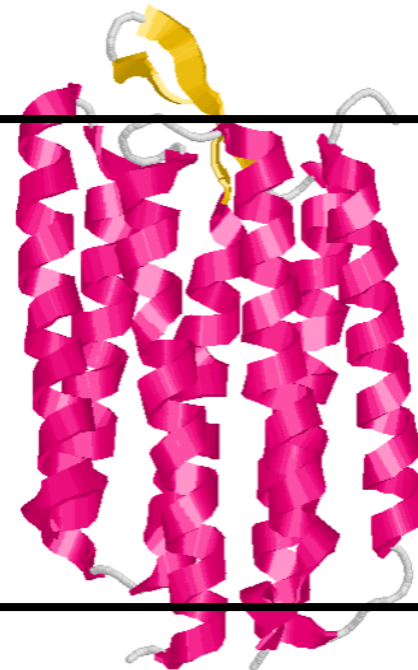


Porin  
(*Rhodobacter capsulatus*)

Inner Membrane

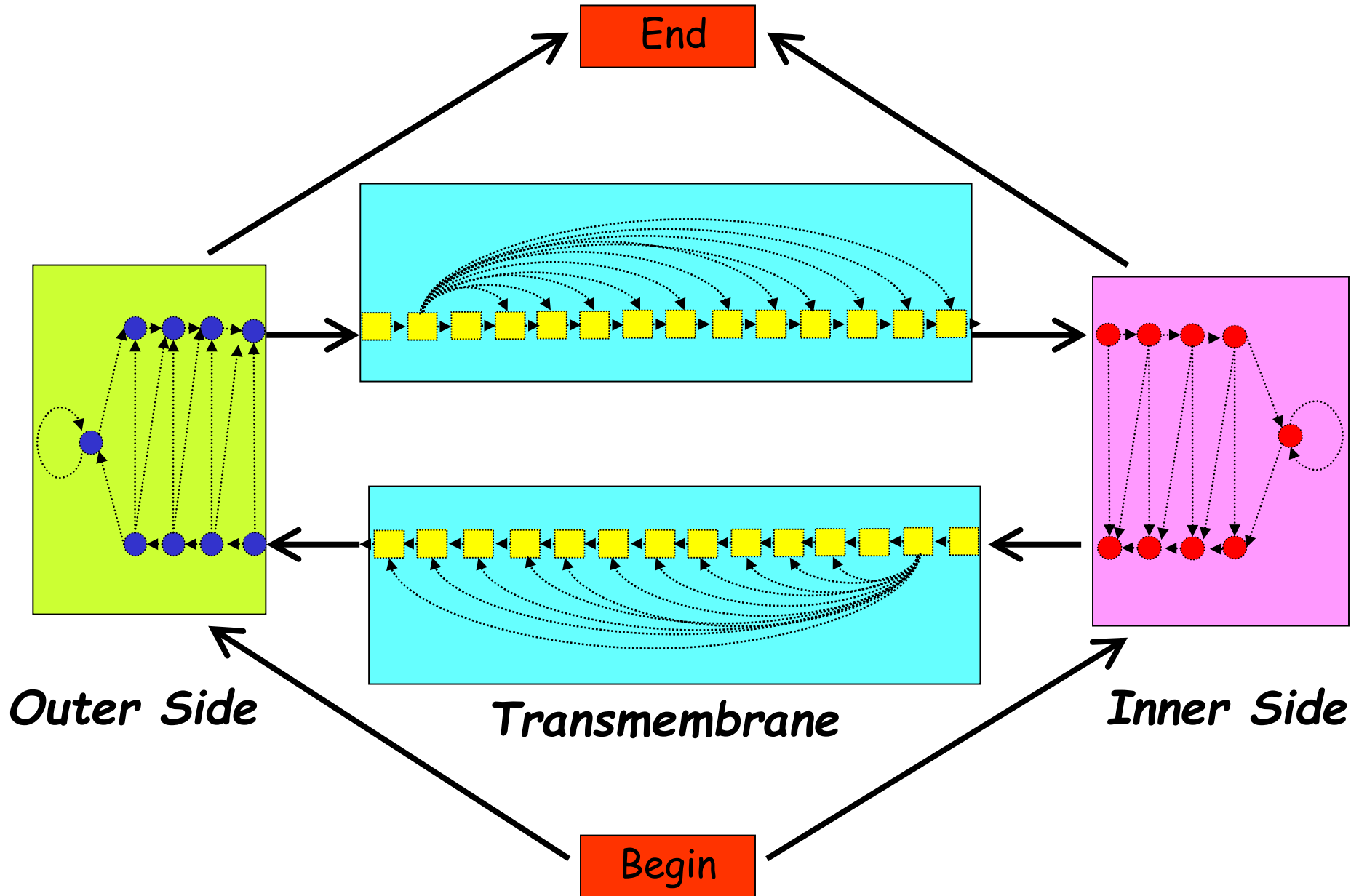
$\alpha$ -helices

Bilayer

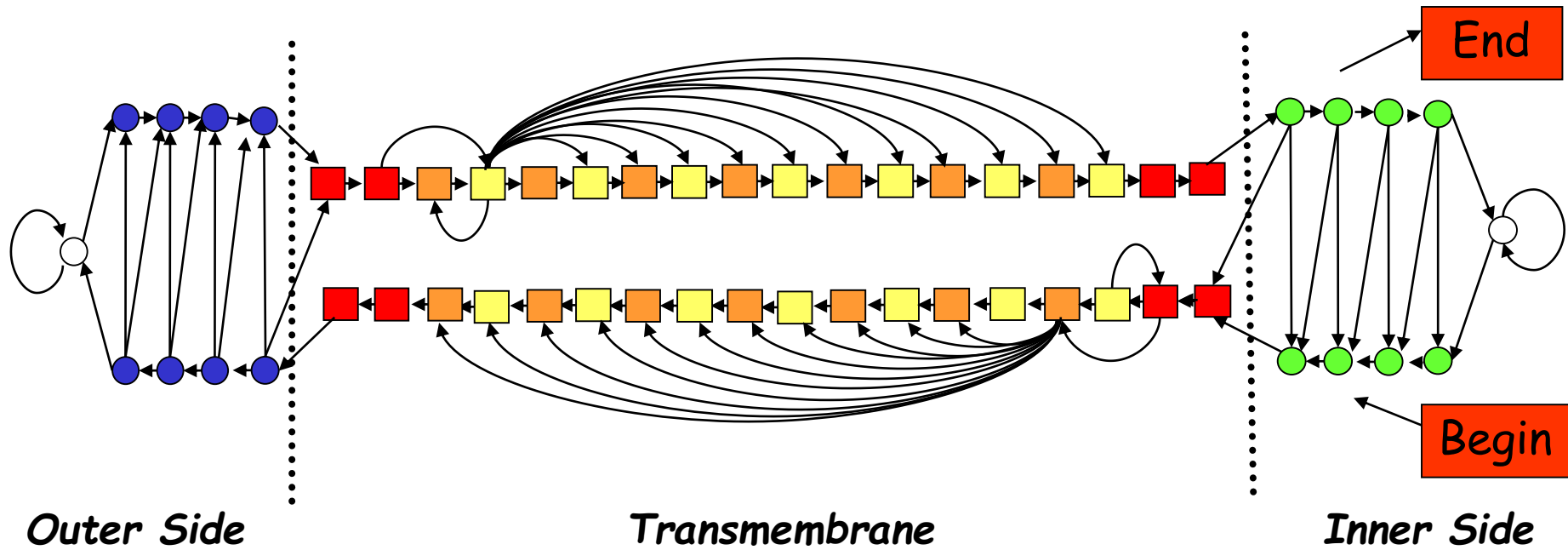


Bacteriorhodopsin  
(*Halobacterium salinarum*)

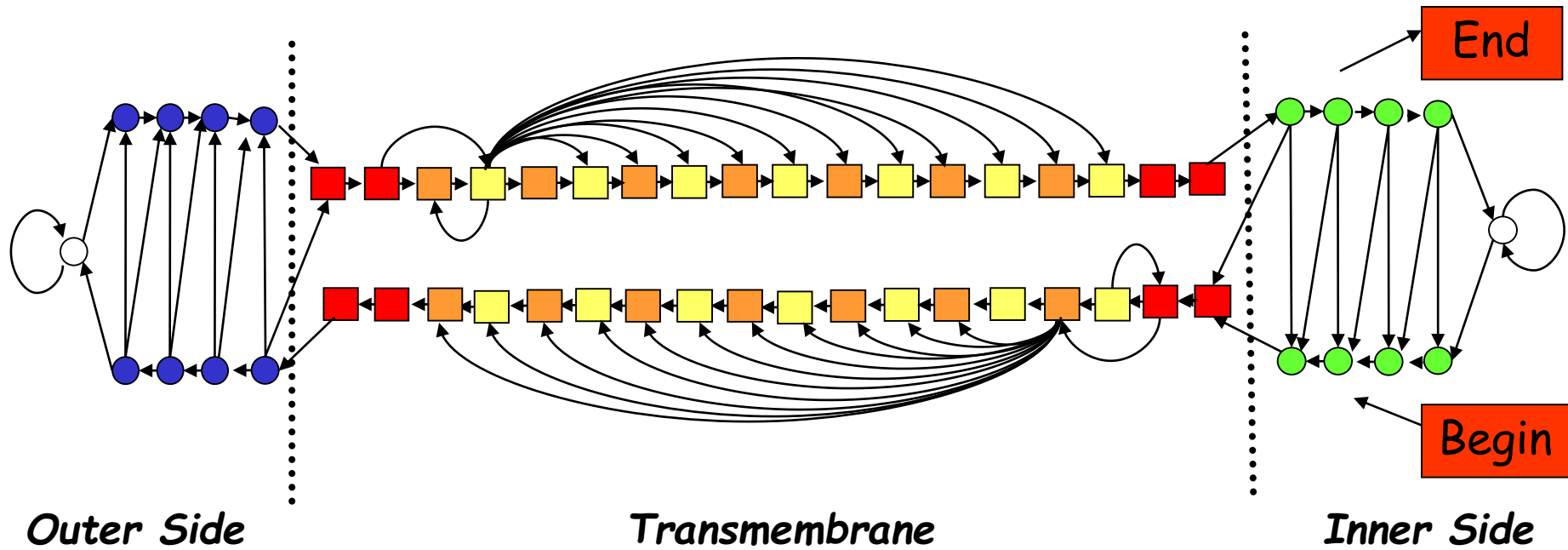
# A generic model for membrane proteins (TMHMM)



# Model of $\beta$ -barrel membrane proteins



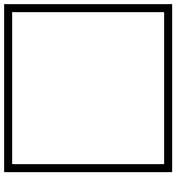
# Model of $\beta$ -barrel membrane proteins



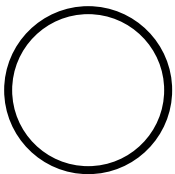
*Outer Side*  
**Labels:**

*Transmembrane*

*Inner Side*



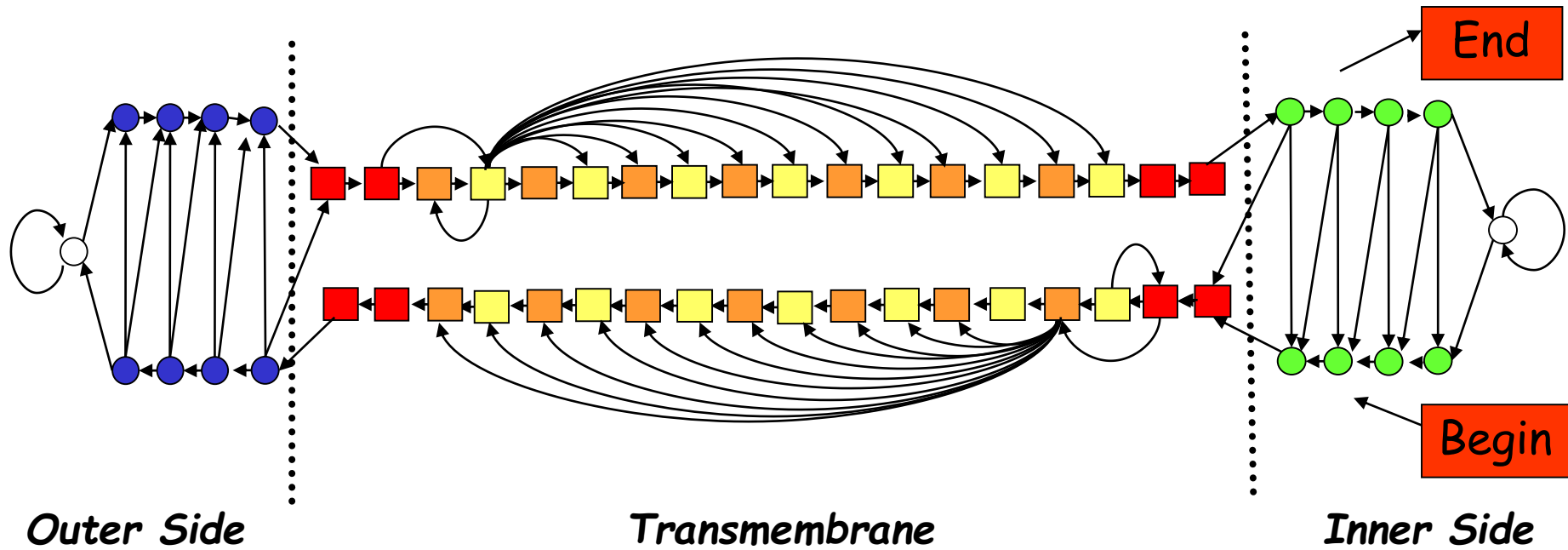
Transmembrane states



Loop states

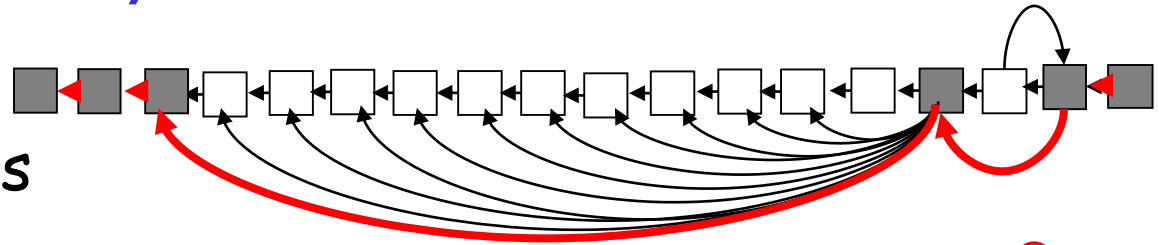


# Model of $\beta$ -barrel membrane proteins

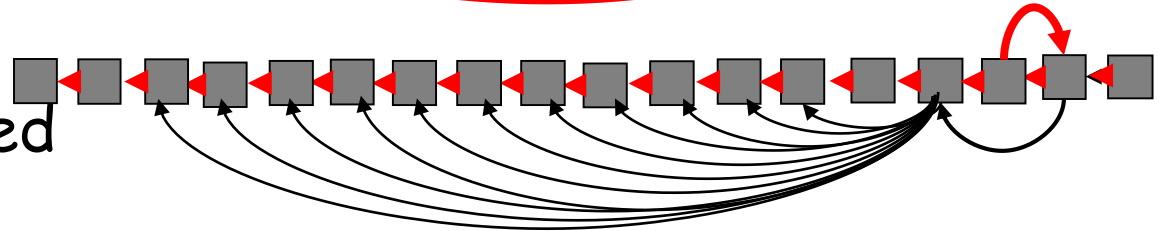


*Length of transmembrane  $\beta$ -strands:*

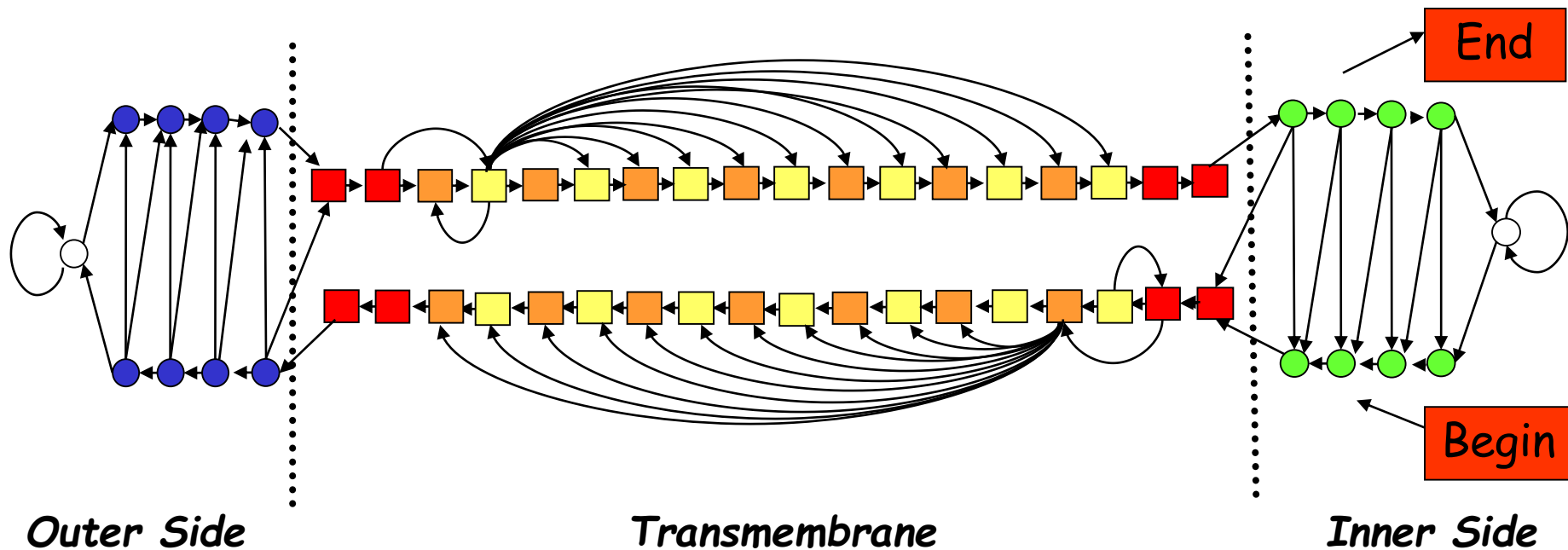
Minimum: 6 residues



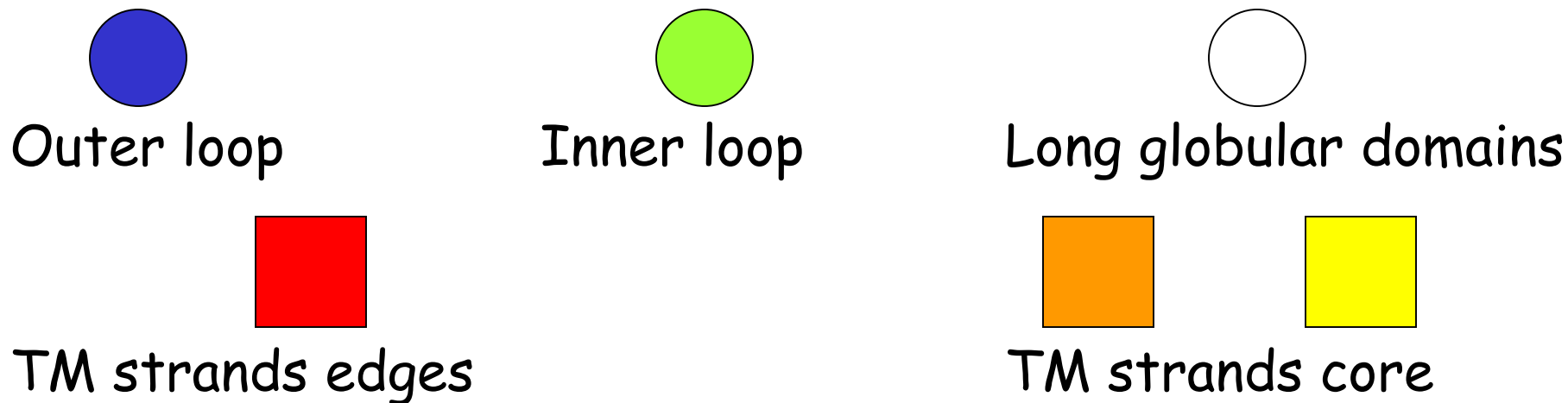
Maximum: unbounded



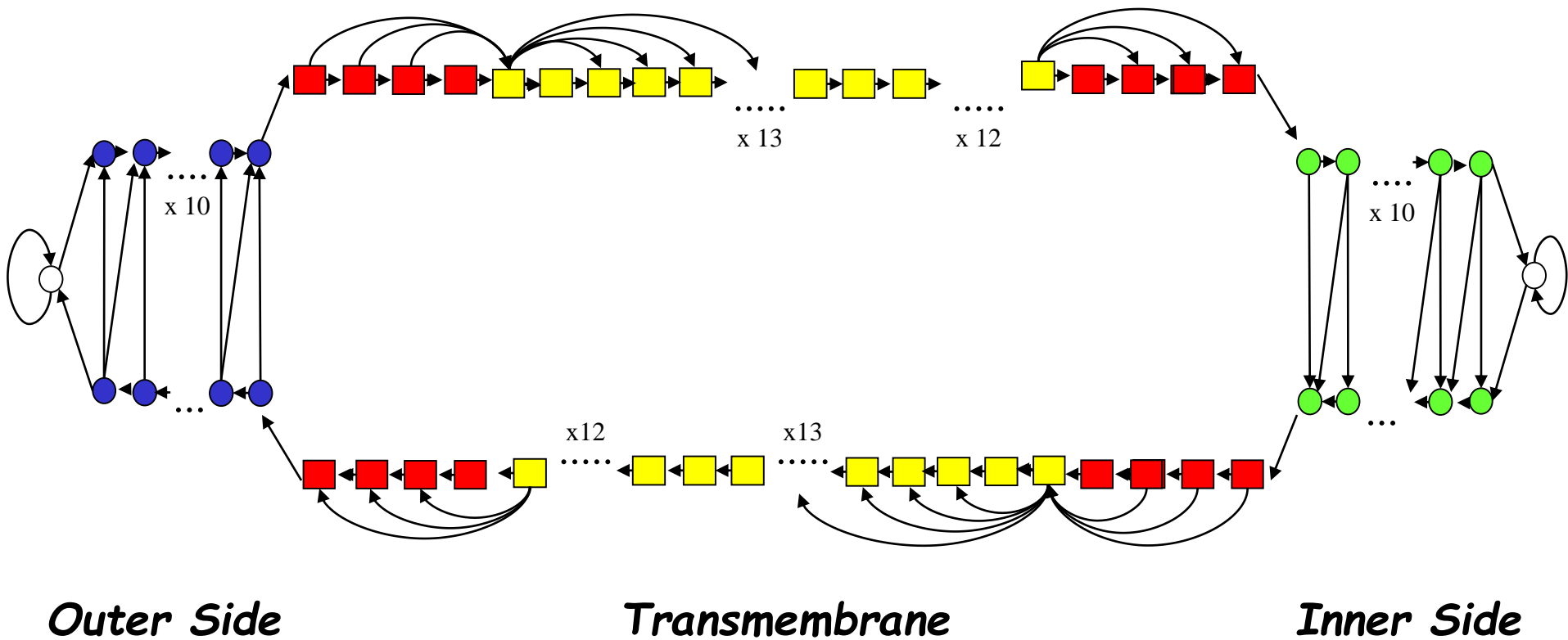
# Model of $\beta$ -barrel membrane proteins



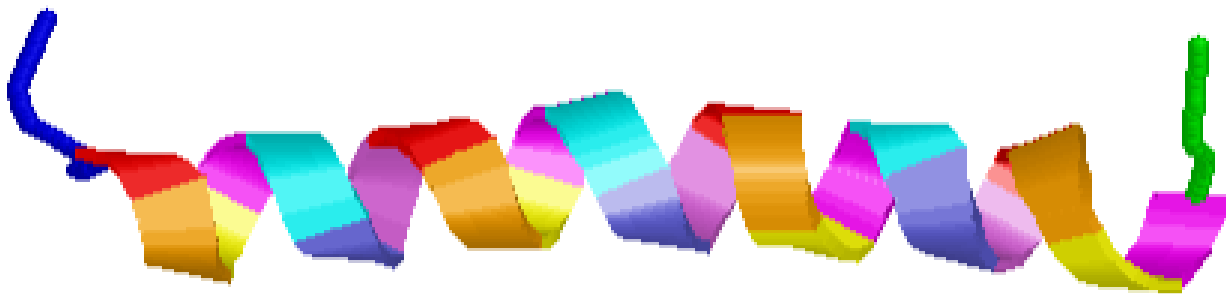
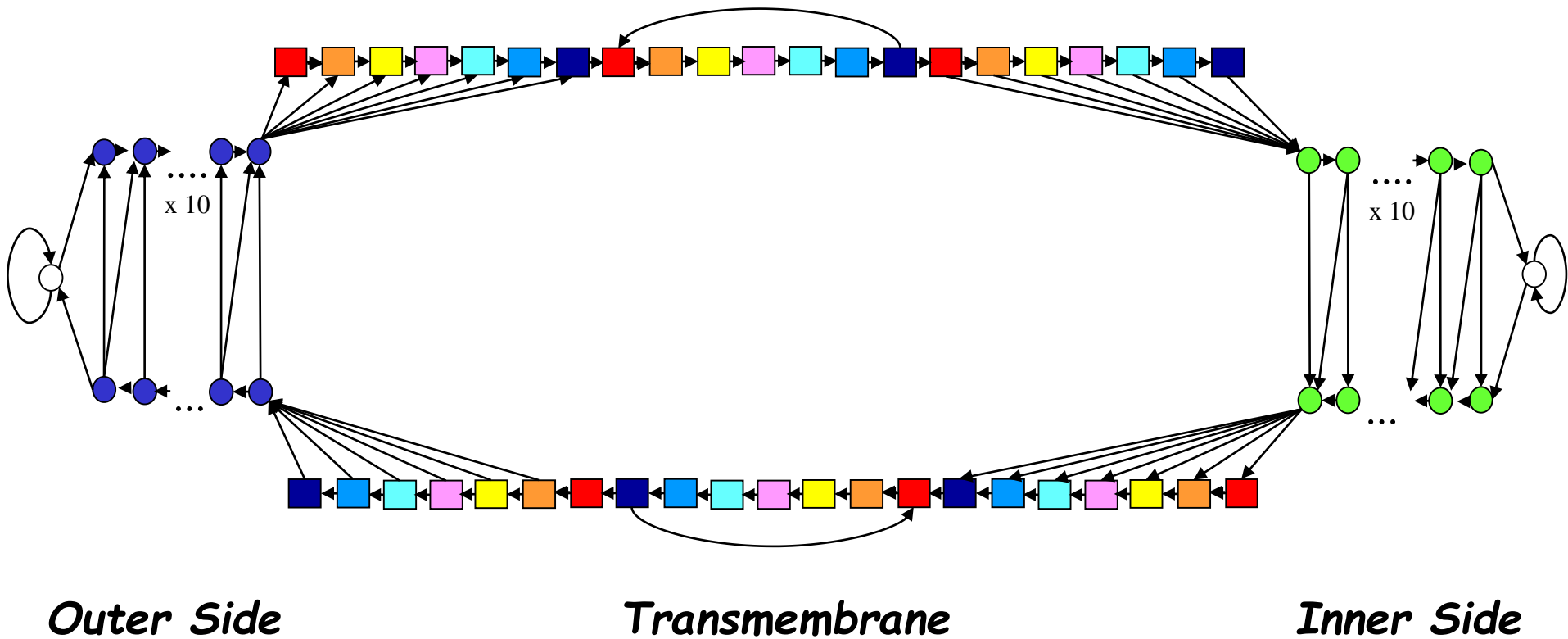
**Six different sets of emission parameters:**



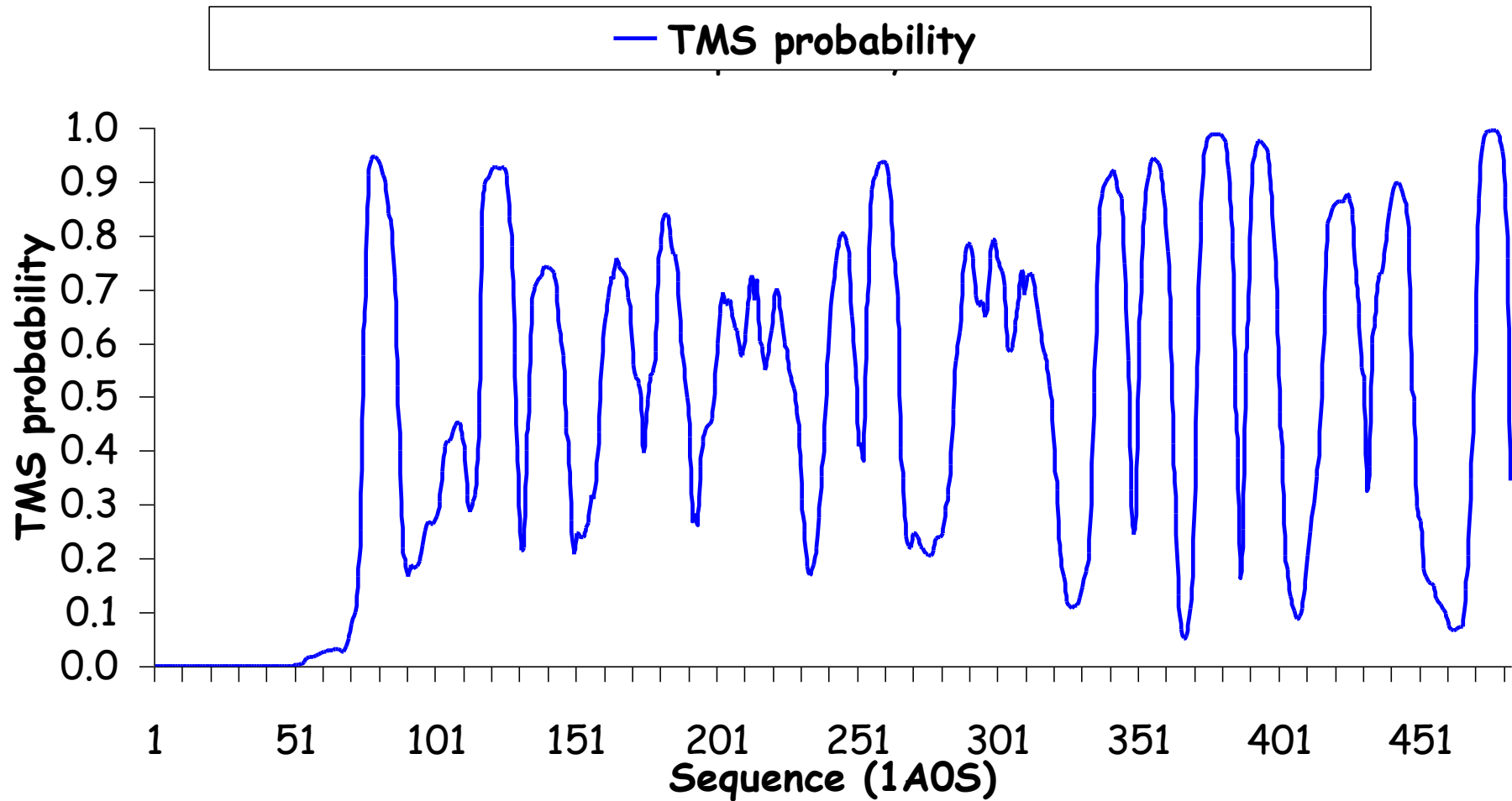
# Model of $\alpha$ -helix membrane proteins (HMM1)



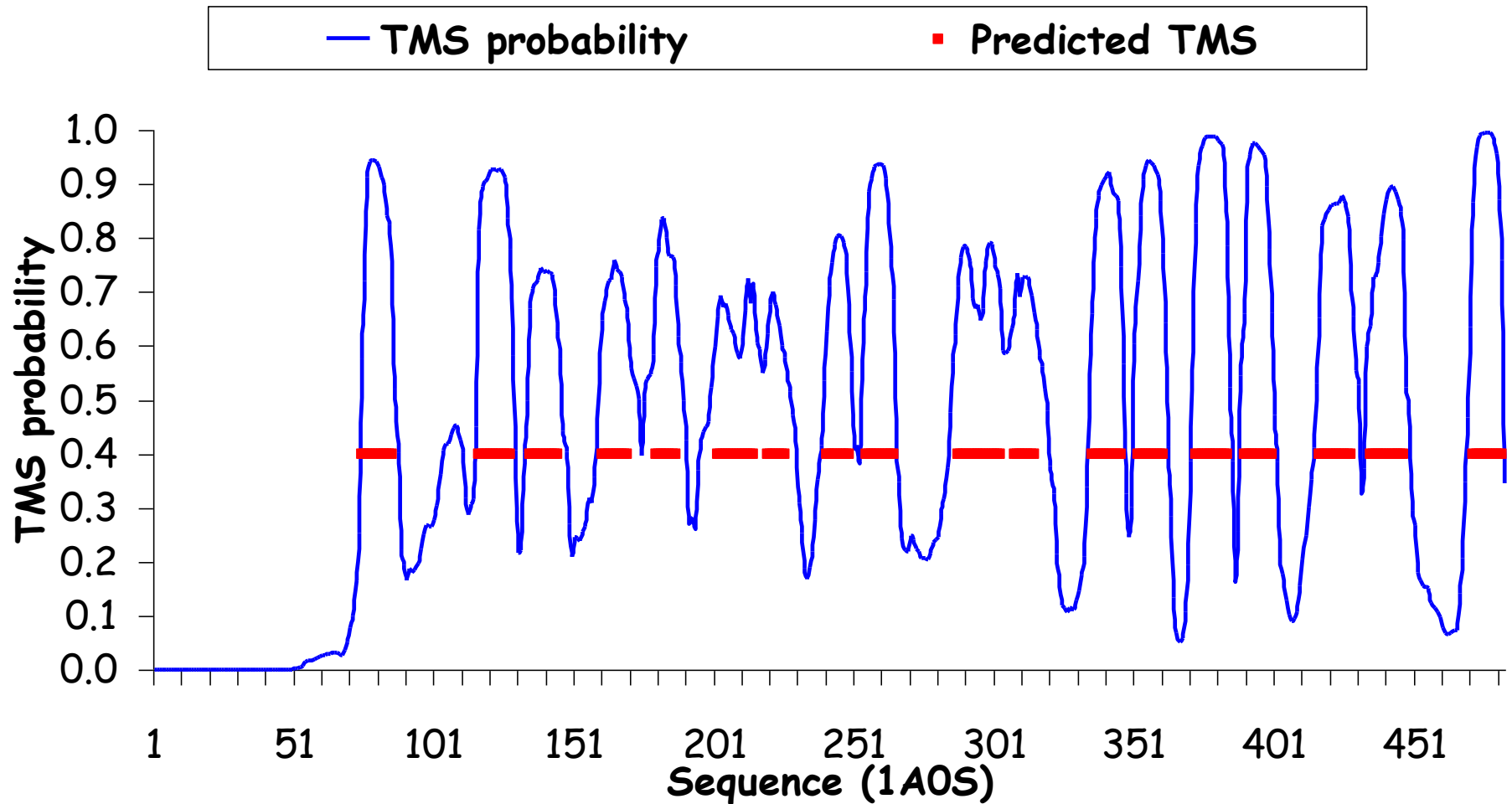
# Model of $\alpha$ -helix membrane proteins (HMM2)



# Dynamic programming filtering procedure

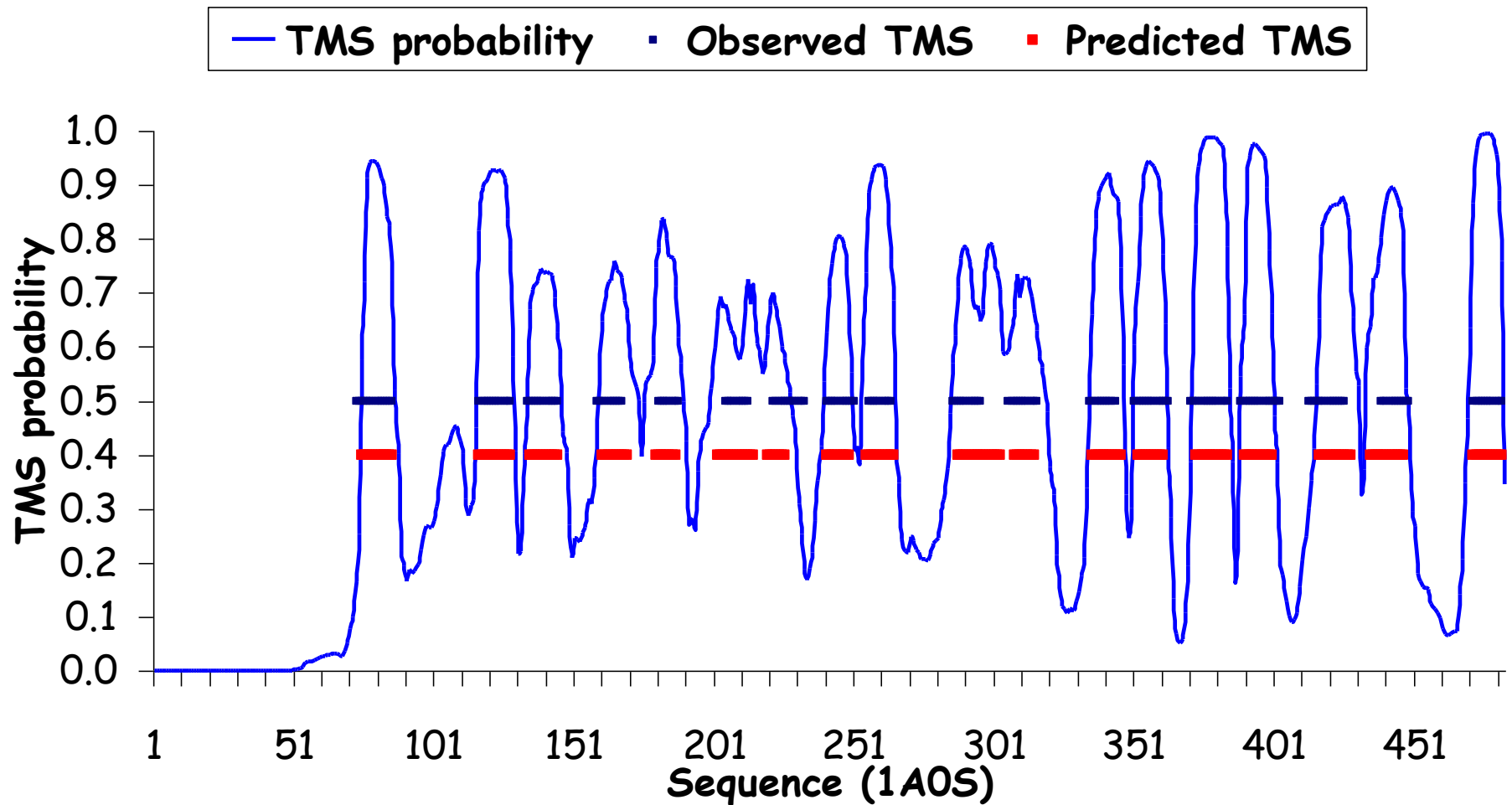


# Dynamic programming filtering procedure



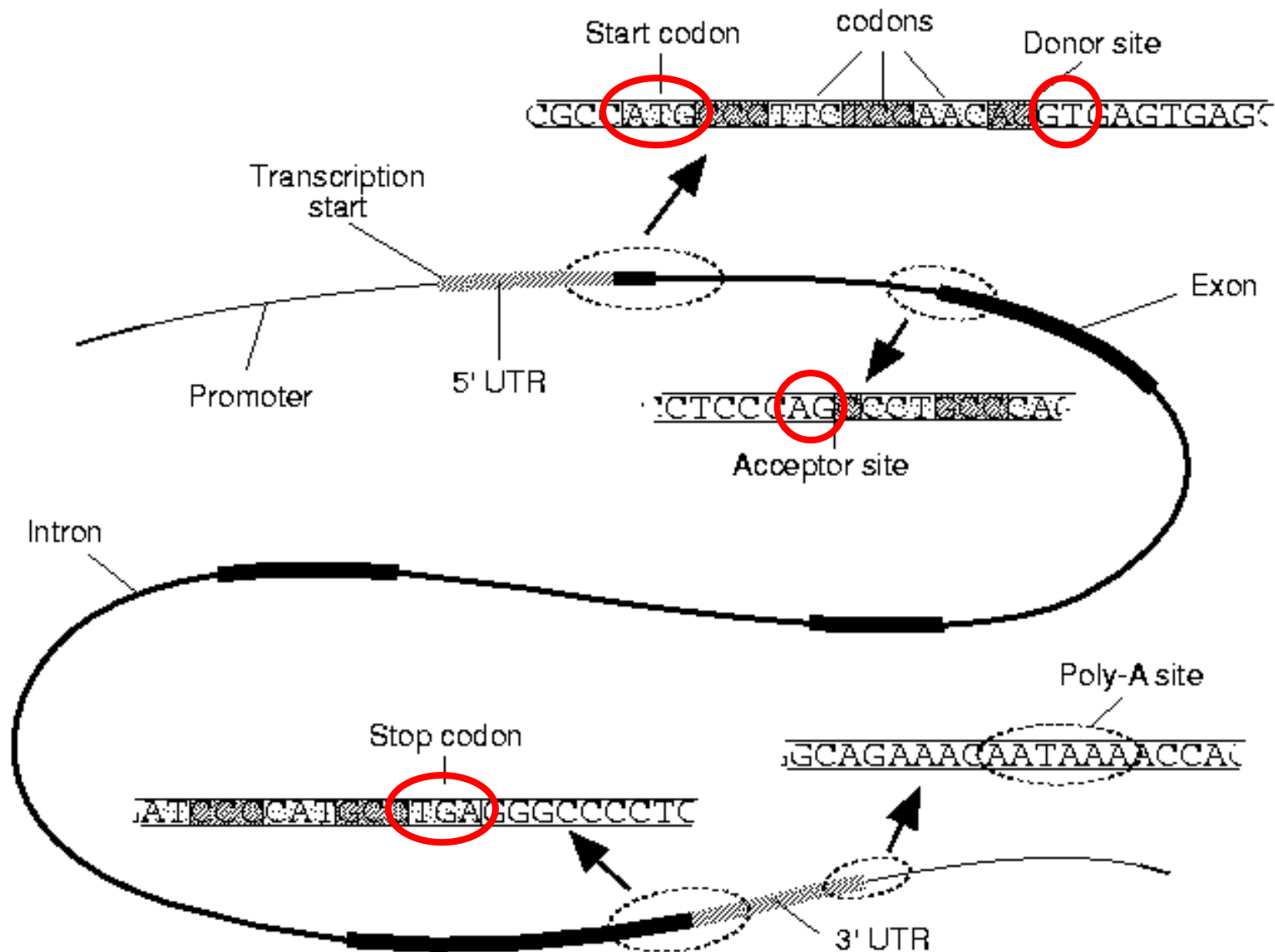
Maximum-scoring subsequences with constrained segment length  
and number

# Dynamic programming filtering procedure



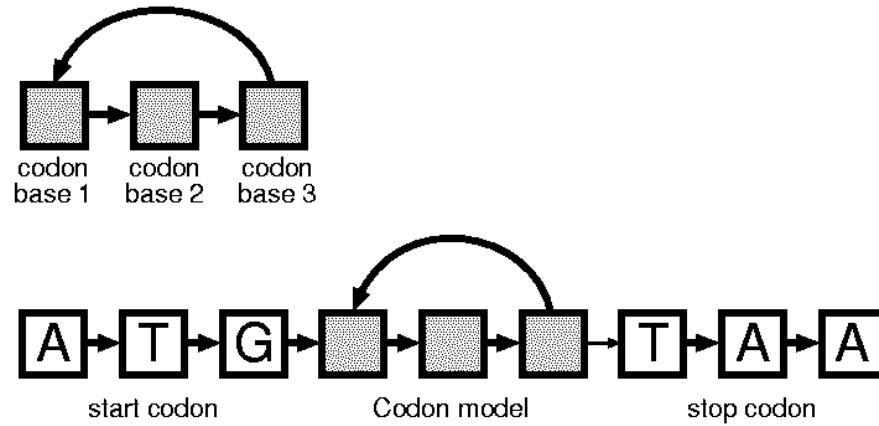
Maximum-scoring subsequences with constrained segment length  
and number

# Eukaryotic gene structure

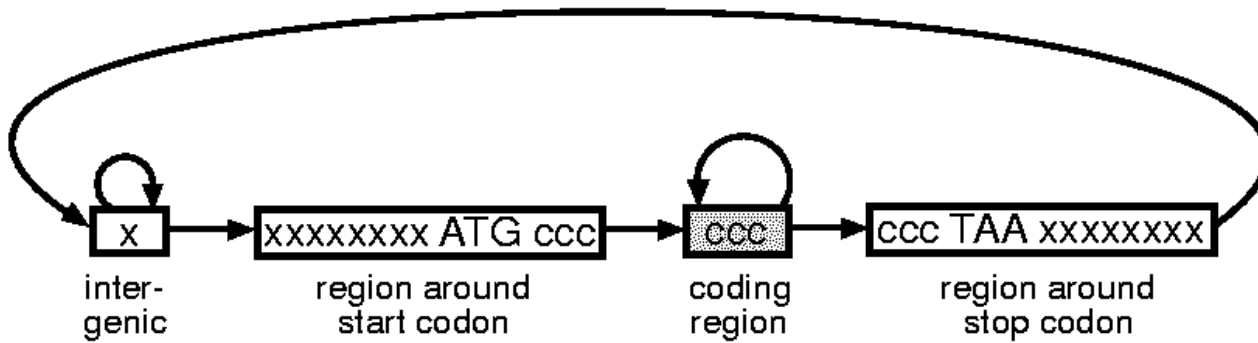




# Simple model for coding regions



# Simple model for unspliced gene



# Simple model for spliced gene

