



# **BIO-463**

## **Genomics and bioinformatics**

### **Lecture 3: Hidden Markov Models**

**Professors: Jacques Rougemont, Anne-Florence Bitbol, Raphaëlle Luisier**

# **EPFL**

# Probability: notation and formulas

$P(x)$  probability of event  $x$

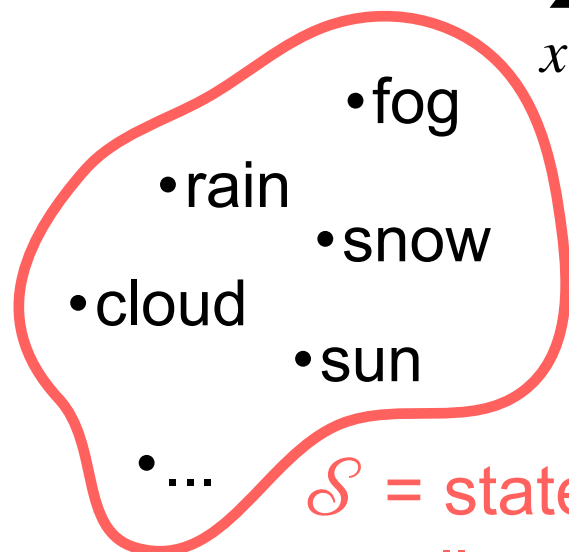
$P(x | y)$  probability of  $x$  knowing  $y$  (**conditioned on**  $y$ )

$P(x, y)$  probability of  $x$  and  $y$

$$P(x, y) = P(x | y)P(y)$$

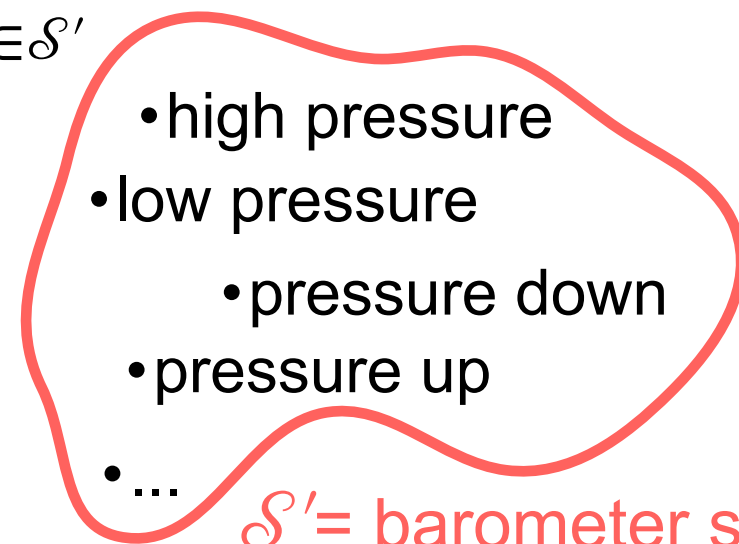
If  $x, y$  independent:  $P(x, y) = P(x)P(y)$ , therefore  $P(x | y) = P(x)$

$$\sum_{x \in \mathcal{S}} P(x) = 1$$



$\mathcal{S}$  = state space  
= all weather conditions

$$\sum_{y \in \mathcal{S}'} P(x | y)P(y) = P(x)$$



$\mathcal{S}'$  = barometer state

# Bayes "Theorem"

$D$  = "data"

$\theta$  = "parameters"

$$\underbrace{P(\theta | D)}_{\text{posterior}} = \frac{\underbrace{P(D | \theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}}{P(D)} = \frac{P(D | \theta)P(\theta)}{\sum_y P(D | y)P(y)}$$

Rare disease: 1/10'000

disease state

	d (yes)	h (no)	
+	10(-0.01)	1000	1010
-	0(+0.01)	98990	98990
	10	99990	100000

Diagnostic test:

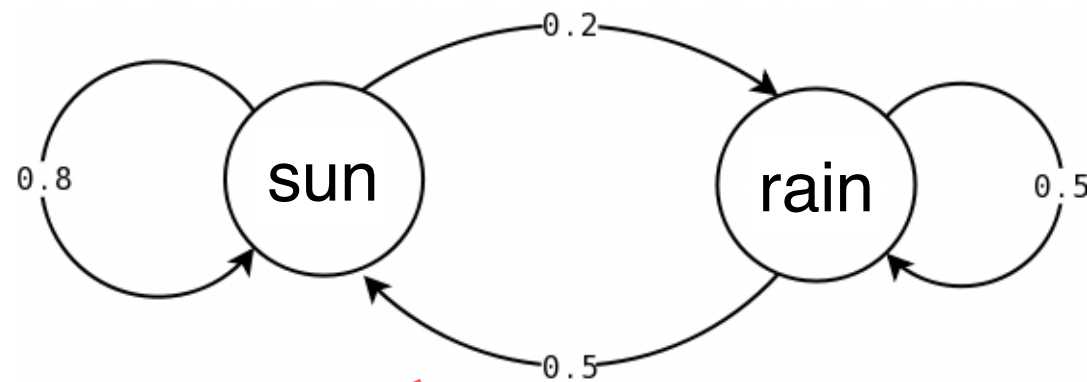
FP (false positive) = 1/100

Sensitivity = 99.9% (FN = 1/1000)

$$\begin{aligned} P(d|+) &= \frac{P(+|d)P(d)}{P(+|d)P(d) + P(+|h)P(h)} \\ &= \frac{0.999 \cdot 10^{-4}}{0.999 \cdot 10^{-4} + 0.01 \cdot (1 - 10^{-4})} \\ &\approx \frac{10^{-4}}{10^{-4} + 10^{-2}} = \frac{1}{101} \end{aligned}$$

# Markov models

- discrete time evolution
- finite state space
- fixed set of transition probabilities



$$M = \begin{pmatrix} \text{sun} & \text{rain} \\ 0.8 & 0.5 \\ 0.2 & 0.5 \end{pmatrix} \begin{matrix} \text{today} & \text{tomorrow} \\ \text{sun} \\ \text{rain} \end{matrix}$$

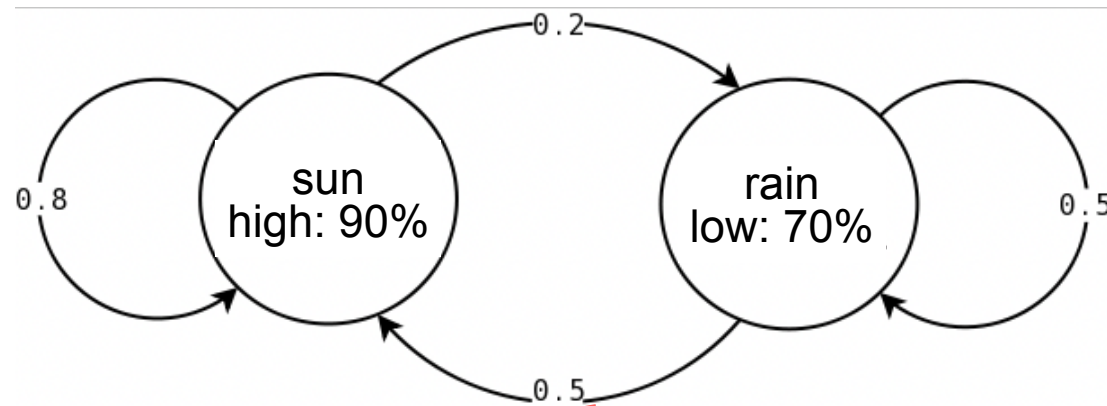
$P(S_{n+1} = \mathbf{r} \mid S_n = \mathbf{s})$

model generates  
plausible sequences

ssssrrssssrrssssrrssssrrssssrrssssrrssssrrssssrr...

$$\begin{pmatrix} p_s(t+1) \\ p_r(t+1) \end{pmatrix} = M \cdot \begin{pmatrix} p_s(t) \\ p_r(t) \end{pmatrix}$$
$$P(S_{n+1} = \mathbf{s}) = \sum_{\sigma} P(S_{n+1} = \mathbf{s} \mid S_n = \sigma) P(S_n = \sigma)$$

# Hidden Markov models (HMM)



transition matrix

$$M = \begin{pmatrix} \text{sun} & \text{rain} \\ 0.8 & 0.5 \\ 0.2 & 0.5 \end{pmatrix} \begin{matrix} \text{sun} \\ \text{rain} \end{matrix}$$

emission matrix

$$E = \begin{pmatrix} \text{sun} & \text{rain} \\ 0.9 & 0.3 \\ \textcircled{0.1} & 0.7 \end{pmatrix} \begin{matrix} \text{high} \\ \text{low} \end{matrix}$$

hidden states

emitted symbols

$$P(T_n = \mathbf{l} | S_n = \mathbf{s})$$

HMM generates

Observed sequence

Hidden sequence

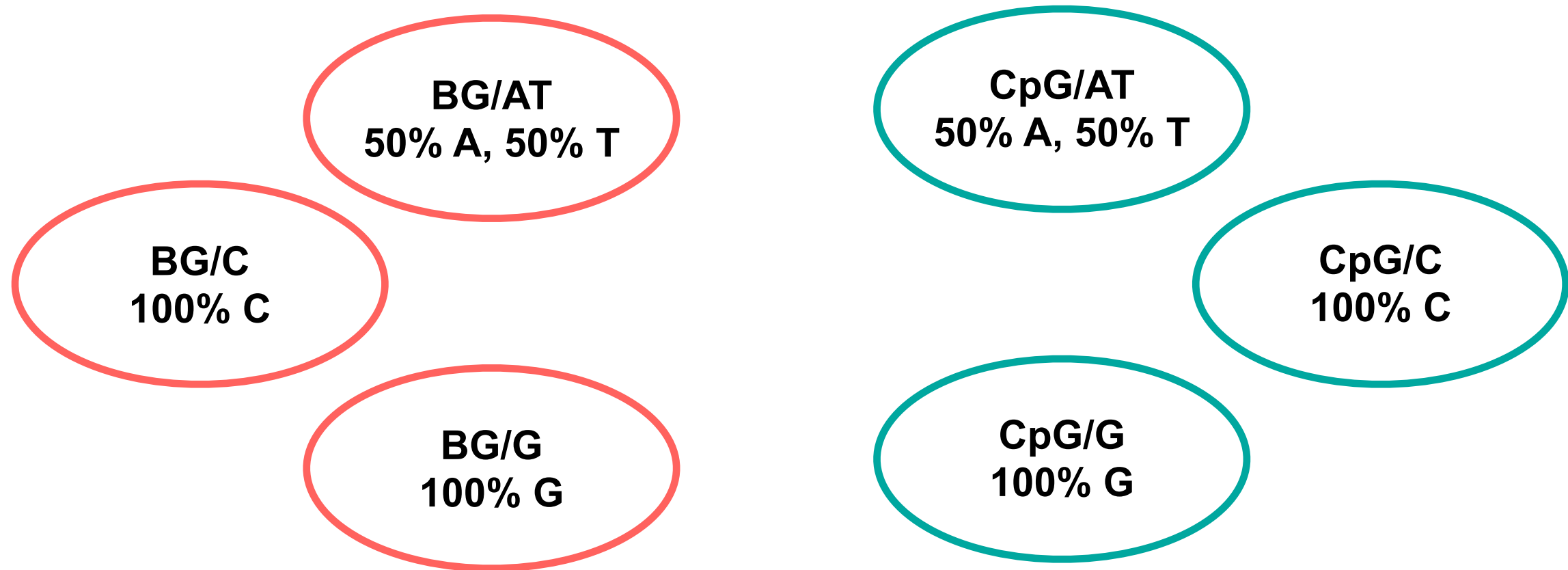
**T** = hhlhhhl111hhlhhlhhlh111hh11111h1h111h11h

**S** = ssssrssssrsssrssssssrrrrsssssssrsssssr

time →

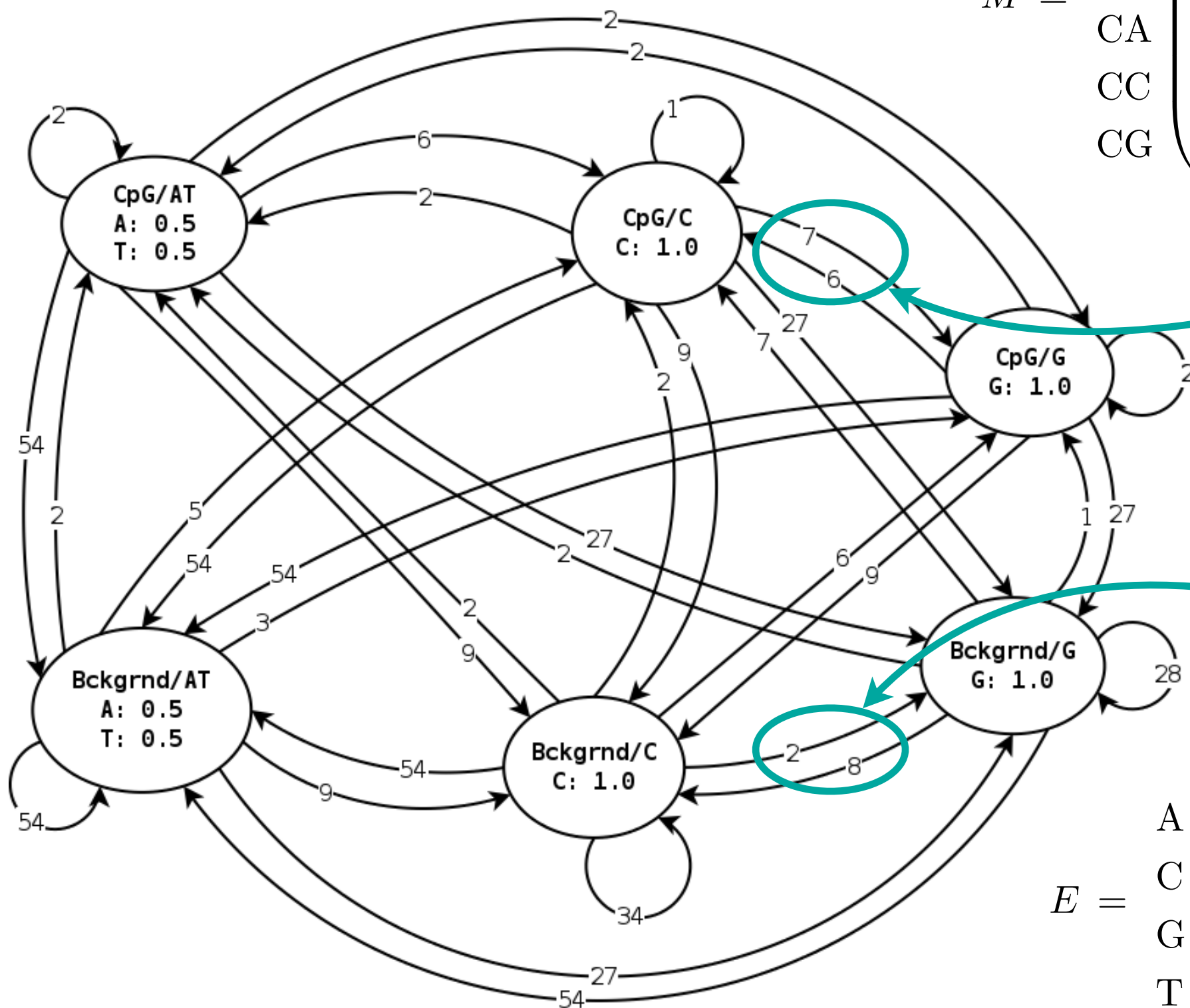
# Example: CpG islands

- Human genome: 30% A/T, 20% C/G
  - Frequency of GC ~ 4%, CC and GG ~ 5%, but CG < 1% !
  - However there are regions called CpG islands containing many consecutive CG
- 
- HMM states: Background (not a CpG island) or CpG
  - States are further divided into A/T, C, G



# Example: CpG islands

$$M = \begin{matrix} & \begin{matrix} BA & BC & BG & CA & CC & CG \end{matrix} \\ \begin{matrix} BA \\ BC \\ BG \\ CA \\ CC \\ CG \end{matrix} & \begin{pmatrix} .54 & .54 & .54 & .54 & .54 & .54 \\ .09 & .34 & .08 & .09 & .09 & .09 \\ .27 & .02 & .28 & .27 & .27 & .27 \\ .02 & .02 & .02 & .06 & .02 & .02 \\ .05 & .02 & .07 & .02 & .01 & .06 \\ .03 & .06 & .01 & .02 & .07 & .02 \end{pmatrix} \end{matrix}$$



**CG, GC > CC, GG**

**CG < CC, GC, GG**

$$E = \begin{matrix} & \begin{matrix} BA & BC & BG & CA & CC & CG \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} 0.5 & 0 & 0 & 0.5 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \\ 0.5 & 0 & 0 & 0.5 & 0 & 0 \end{pmatrix} \end{matrix}$$

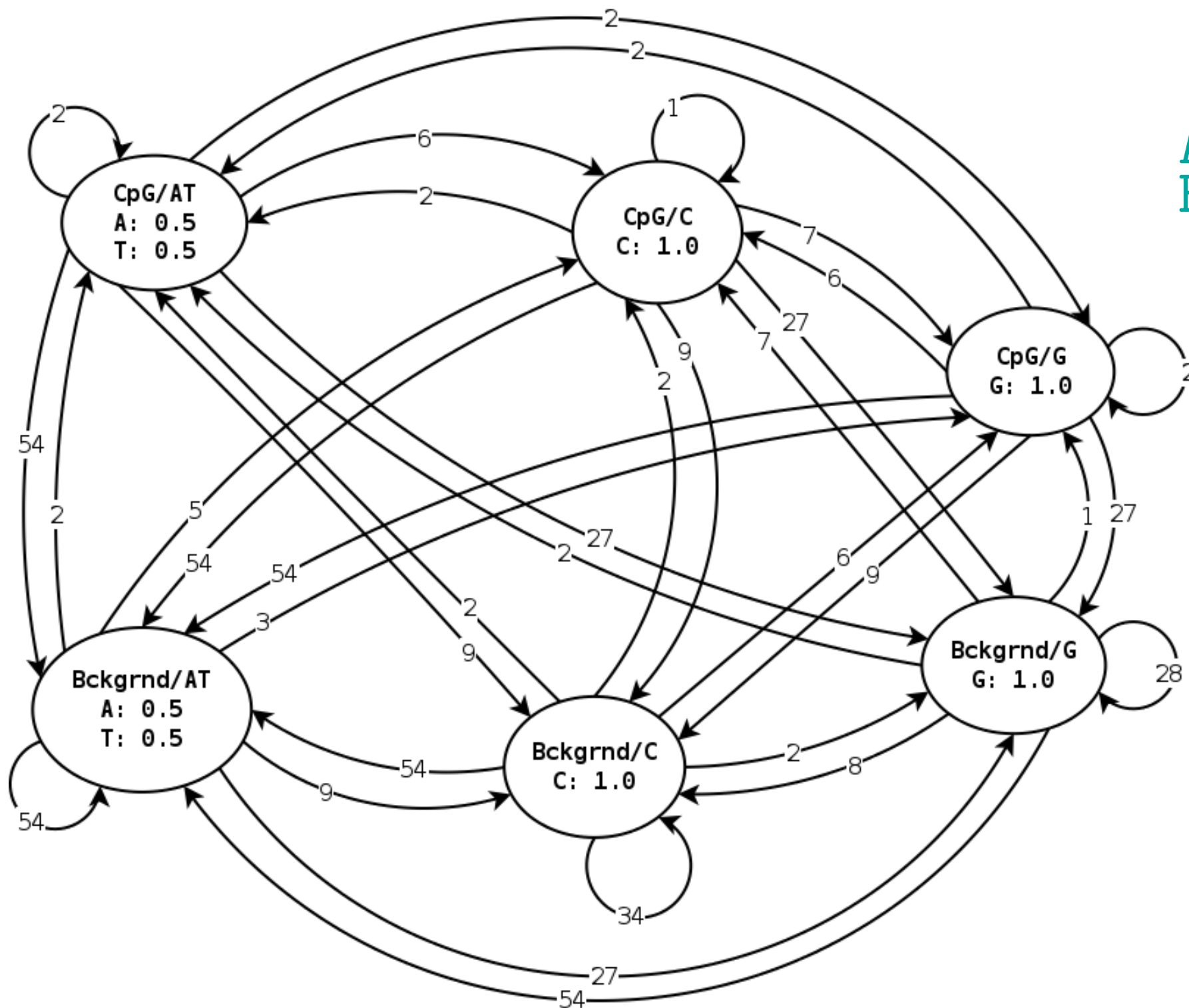
# Example: CpG islands

ACGCGTCGCGA  
BCCCCCCCCCB

$$5\ 7\ 6\ 7\ 2\ 6\ 7\ 6\ 7\ 54 = 2.8e-12$$

ACGCGTCGCGA  
BBBBBBBBBBBB

$$9\ 2\ 8\ 2\ 54\ 9\ 2\ 8\ 2\ 54 = 2.4e-12$$





# Hidden Markov Models (HMM)

## Questions:

- What is the sequence  $\mathbf{S}$  of hidden states most likely to generate the observed symbols  $\mathbf{T}$ ?
- What is the probability of the observed sequence  $\mathbf{T}$ ?
- What are the parameters  $E$ ,  $M$  (emission and transition probabilities) that maximize the probability of  $\mathbf{T}$ ?

## Answers:

- Viterbi algorithm
- Forward / Backward algorithms
- Baum-Welch algorithm

# Hidden Markov Models (HMM)

## Questions:

## Answers:

- What is the sequence  $S$  of hidden states that could have generated the observed symbols  $O$ ?  
- Viterbi algorithm

**Decoding the observations**

- What is the probability of the sequence  $O$  given the model  $\lambda$ ?  
- Forward / Backward algorithms

**Evaluating the model**

- What are the parameters  $E$ ,  $M$  (emission and transition probabilities) that maximize the likelihood of the training data?  
- Welch algorithm

**Optimizing the model**

# HMM formulas

Probability of hidden sequence:

$$P(\mathbf{S}) = P(S_0) \prod_{n=1}^N M(S_n, S_{n-1})$$

initial state

transition

Probability of observed sequence  
(knowing hidden states):

$$P(\mathbf{T}|\mathbf{S}) = \prod_{n=1}^N P(T_n|S_n) = \prod_{n=1}^N E(T_n, S_n)$$

emission from  
hidden state

Bayes theorem:

$$P(\mathbf{S}|\mathbf{T}) = \frac{P(\mathbf{T}|\mathbf{S})P(\mathbf{S})}{P(\mathbf{T})} = \frac{P(\mathbf{T}|\mathbf{S})P(\mathbf{S})}{\sum_{\sigma} P(\mathbf{T}|\sigma)P(\sigma)}$$

Numerical stability:

$$\log P(\mathbf{S}) = \log P(S_0) + \sum_{n=1}^N \log M(S_n, S_{n-1})$$

# Viterbi algorithm

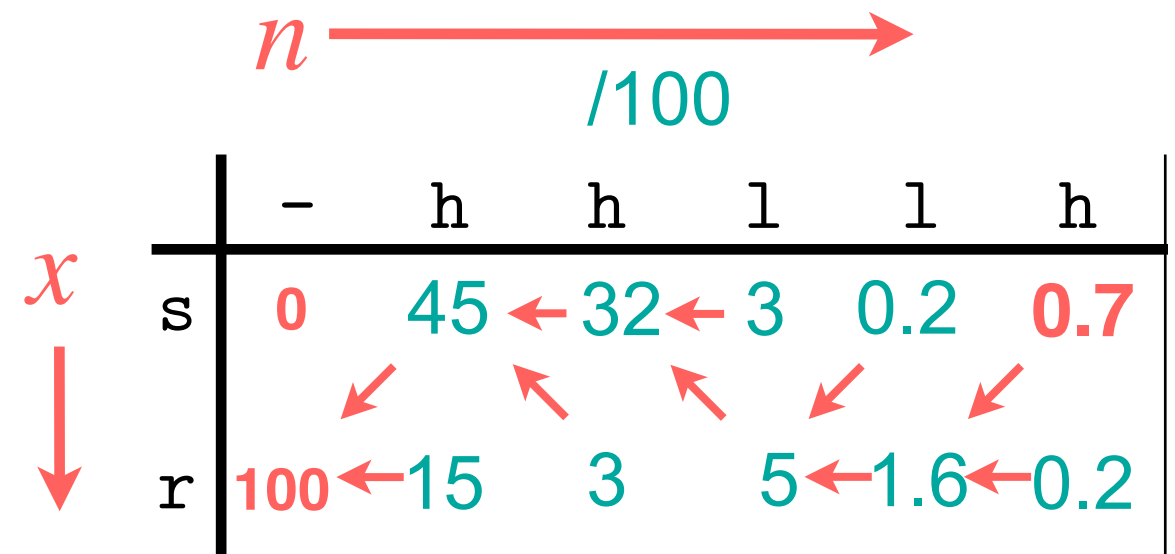
## Example

$\mathbf{T} = \text{hhllh}$

$$V_{n,x} = E(T_n, x) \max_y M(x, y) V_{n-1,y}$$

$$M = \begin{pmatrix} \text{sun} & \text{rain} \\ 0.8 & 0.5 \\ 0.2 & 0.5 \end{pmatrix} \begin{matrix} \text{sun} \\ \text{rain} \end{matrix}$$

$$E = \begin{pmatrix} \text{sun} & \text{rain} \\ 0.9 & 0.3 \\ 0.1 & 0.7 \end{pmatrix} \begin{matrix} \text{high} \\ \text{low} \end{matrix}$$



# Forward / Backward algorithm

- What is the probability of my observation  $\mathbf{T}$ ?  $P(\mathbf{T}) = \sum_{\mathbf{s}} P(\mathbf{T} | \mathbf{s})P(\mathbf{s})$
- Forward score (replace max by sum in Viterbi):

$$F_{n,x} = E(T_n, x) \sum_y M(x, y) F_{n-1,y} = P(T_1, \dots, T_n, S_n = x)$$

- Backward (similar to Forward, move from right to left):

$$B_{n,x} = \sum_y E(T_{n+1}, y) M(y, x) B_{n+1,y} = P(T_{n+1}, \dots, T_N | S_n = x)$$

$$\frac{F_{n,x} \cdot B_{n,x}}{\sum_y F_{N,y}} = \frac{\frac{P(\mathbf{T} | S_n = x) \cdot P(S_n = x)}{P(\mathbf{T})} \cdot P(T_1, \dots, T_n, S_n = x) \cdot P(T_{n+1}, \dots, T_N | S_n = x)}{P(\mathbf{T})} = P(S_n = x | \mathbf{T})$$

$P(\dots T_n, S_n) = P(\dots T_n | S_n)P(S_n)$

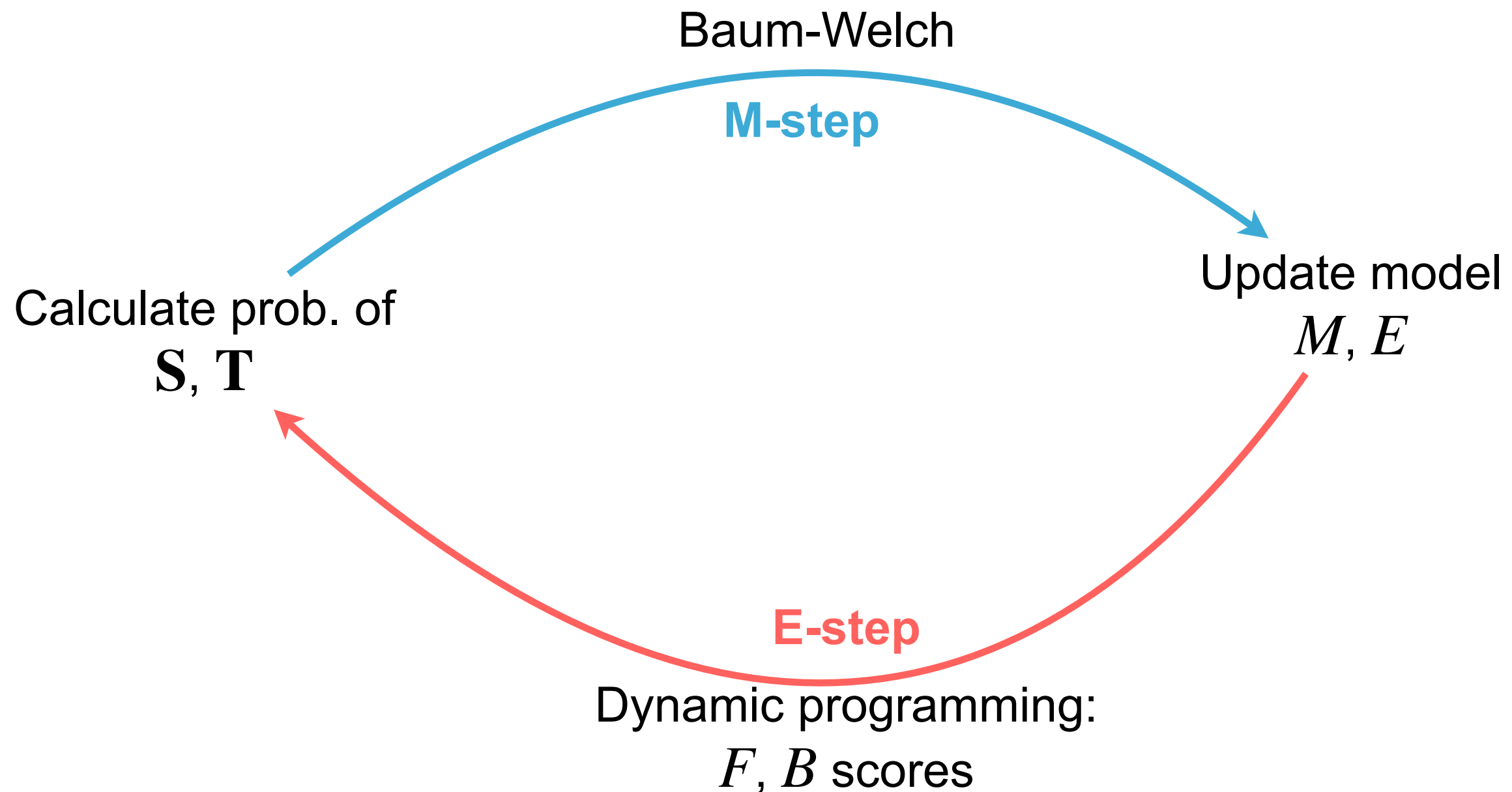
$P(\dots T_n, T_{n+1} \dots | S_n) = P(\dots T_n | S_n)P(T_{n+1} \dots | S_n)$

Bayes

Probability of any hidden state given observed  $\mathbf{T}$

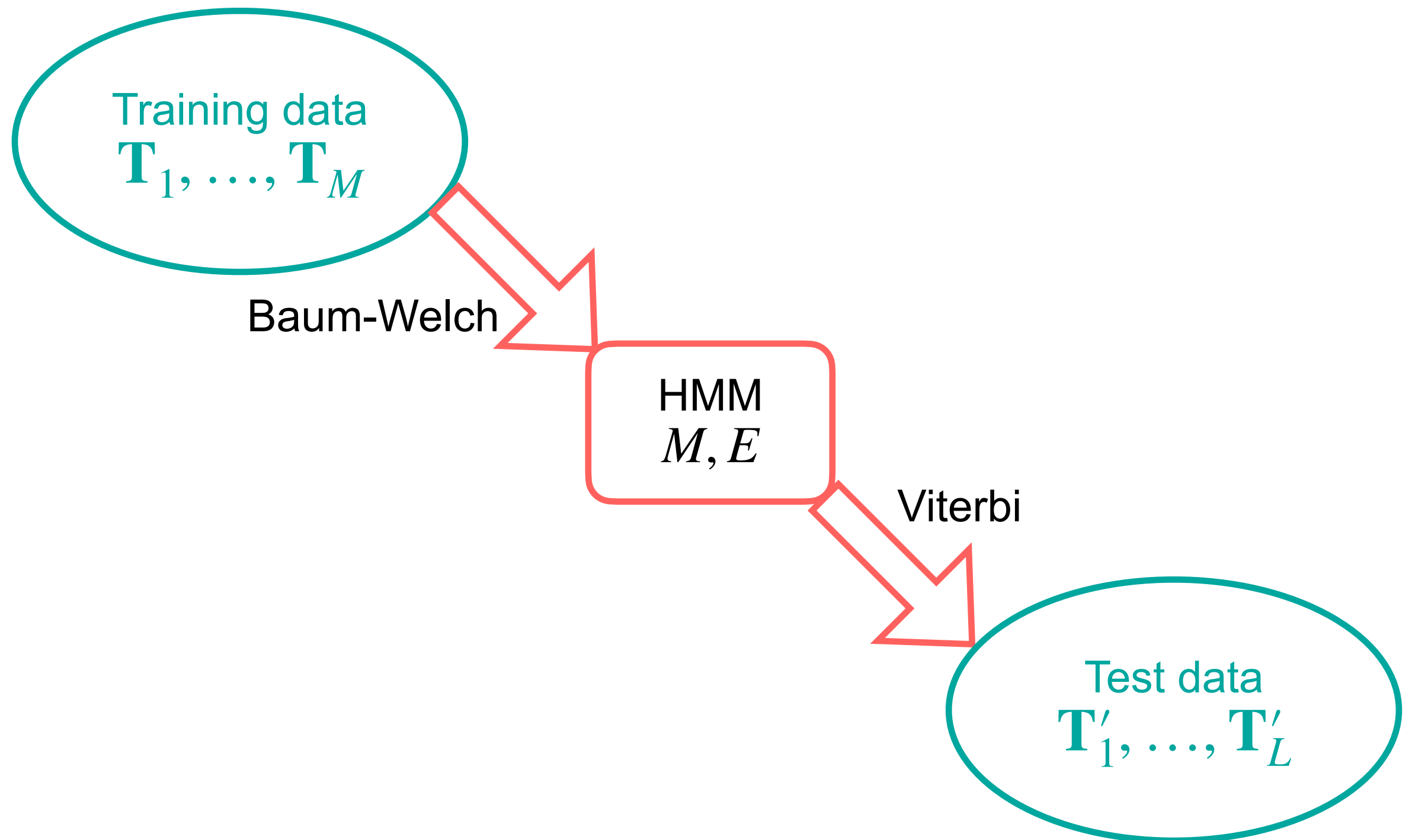
# Baum-Welch optimization

Example of Expectation-Maximization (EM) algorithm

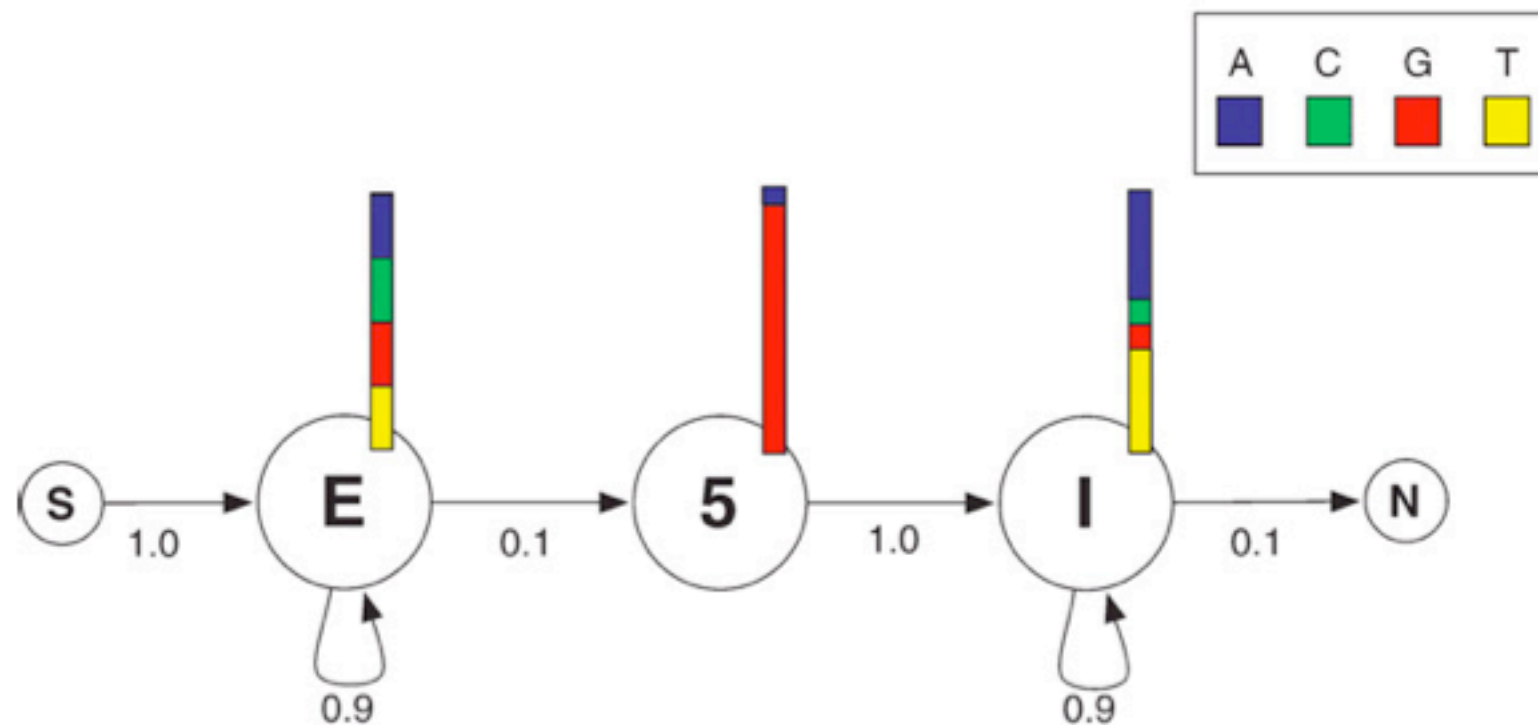


⇒ converges to a model which maximizes  $P(T)$

# Model "learning"



# Modeling spliced genes

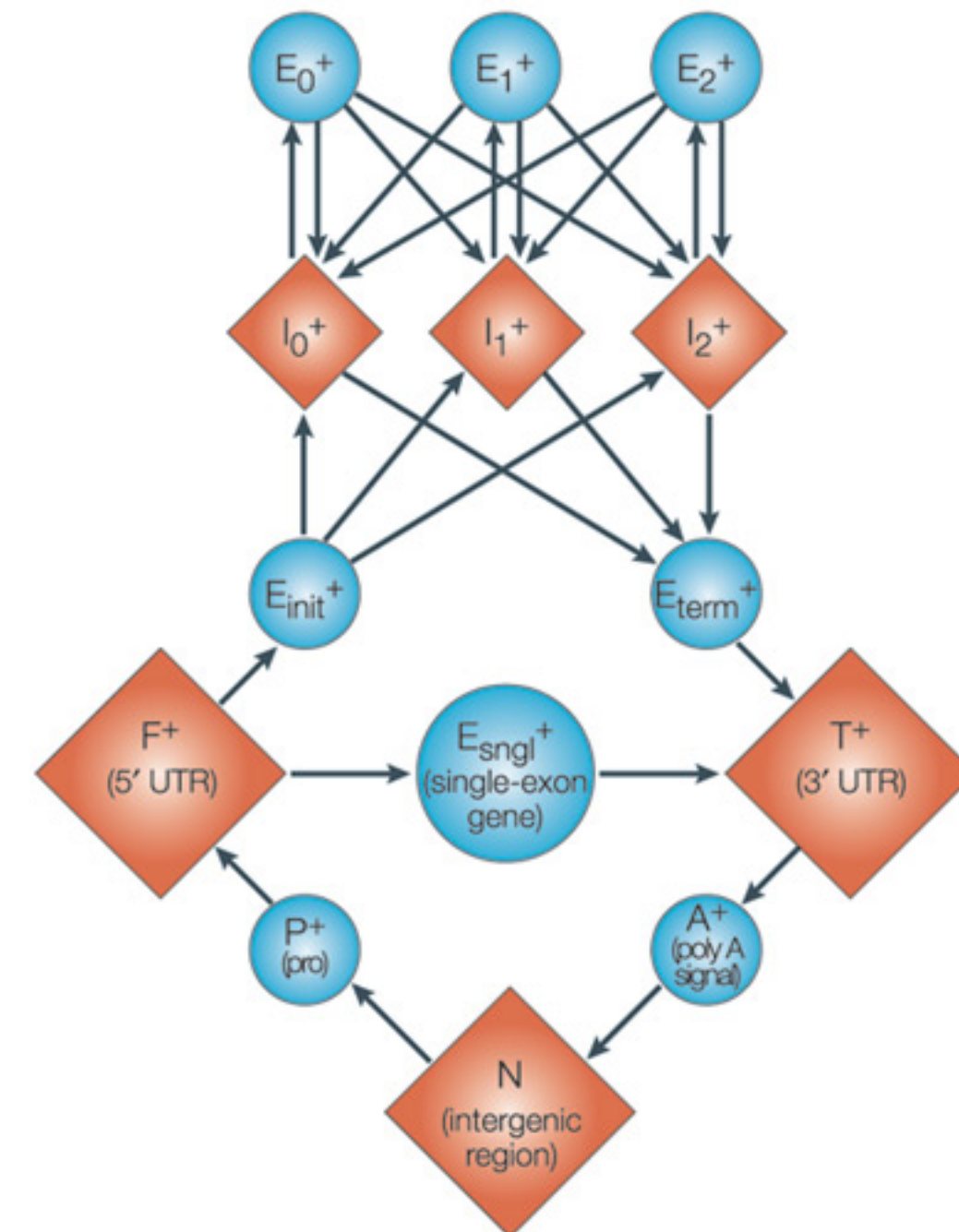
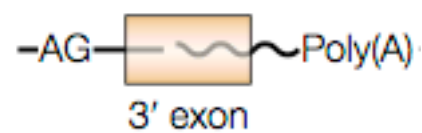
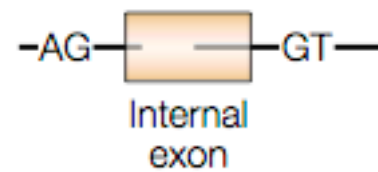
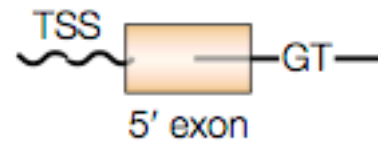


**S** = ... EEEEEEEEE5IIIIIII ...  
**T** = ... CAGTGTAAAGTATCATT ...



# Genscan programm

[wikipedia:GENSCAN](https://en.wikipedia.org/wiki/GENSCAN)



Reverse strand: mirror reflection of above

Nature Reviews | **Genetics**

Zhang MQ, Nat Rev Gen (2002)