

3. Predicting binding sites

- basics of TFBS identification
- defining a background model
- tools
- phylogenetic footprinting
- **including "in-vivo features"**



Institut für Pharmazie und
Molekulare Biotechnologie



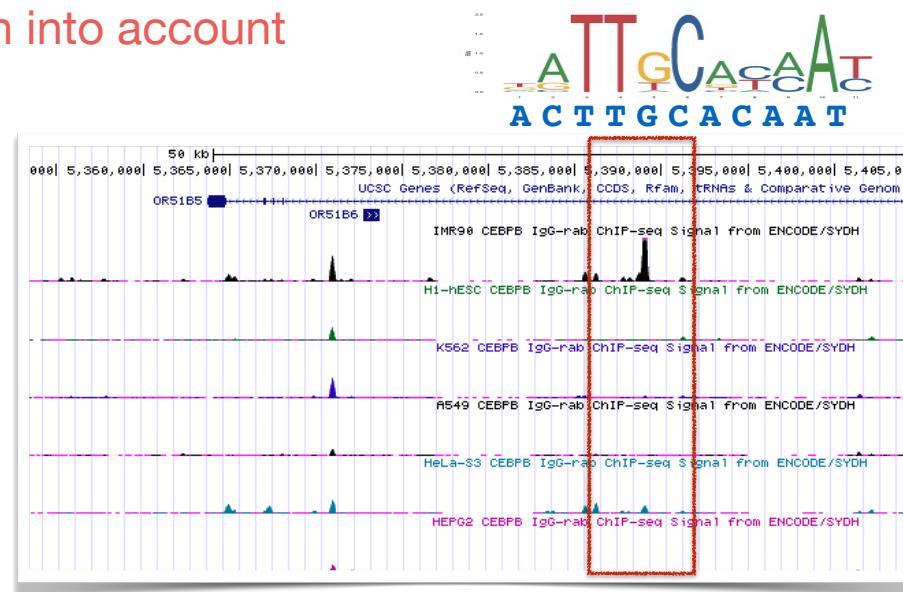
UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Improving TFBS predictions

● Limitations

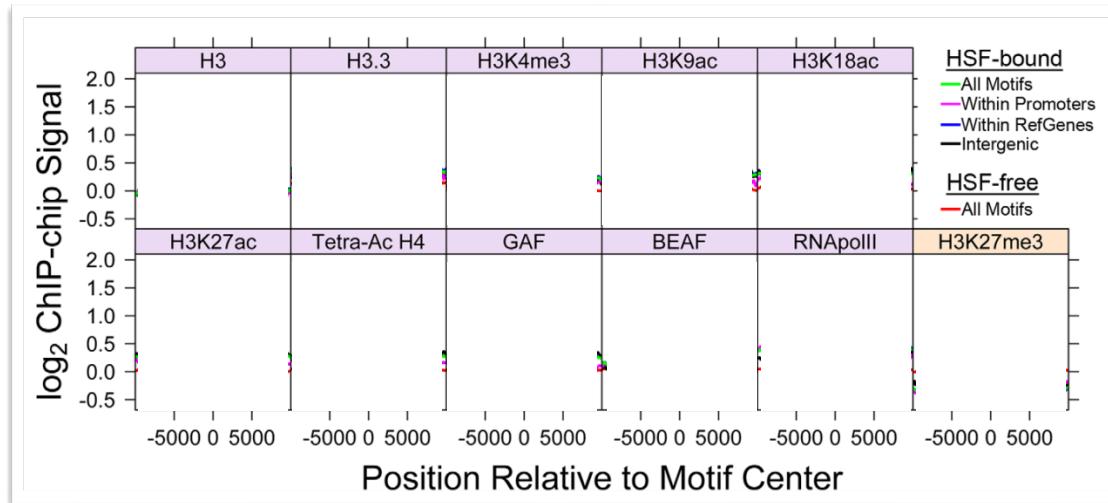
- quality of the matrix (PWMs constructed from few sites are not discriminative, low information content !)
- difficulty to predict **low affinity** binding events
- correct choice of the **background** model
- **in-vivo context** is not taken into account

Why is CEBPB
binding in some
cell-lines and
not in others ?



Motifs are not always binding events

- Compare **in-silico** TFBS to **in-vivo** binding event using ChIP data
- some bona fide motifs are not bound in-vivo : why ?
- example in Drosophila : heat-shock factor (HSF)
 - 464 ChIP peaks containing a HSF-motif ($p < 0.001$)
 - 708 unbound motifs (with $p < 5e-6$)



Bound sites have:

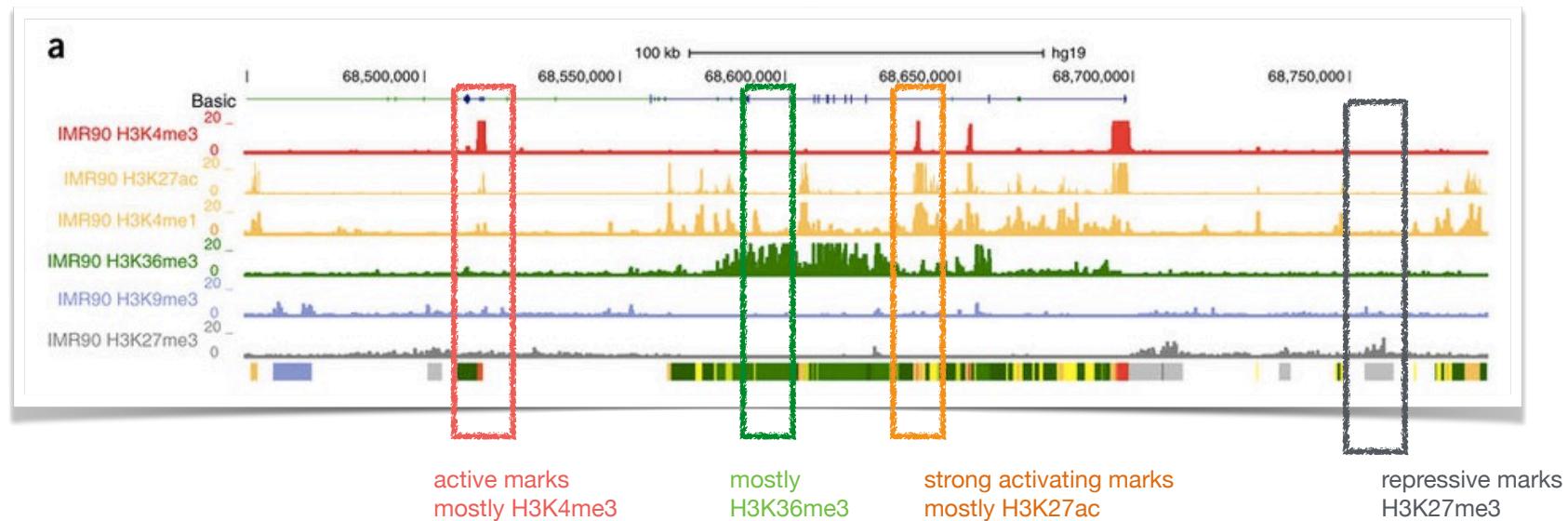
- high levels of lysine acetylation
- high levels of Pol2 binding
- low levels of H3K27me3 (repressive mark related to polycomb repression)

[Guertin et al., PLoS Gen. (2010)]

Histone code

Mark	Interpretation
H3K4me1	activating mark; found at promoters and enhancers
H3K4me3	mark of active and open gene promoters
H3K27ac	mark of active promoters and enhancers
H3K36me3	mark of actively transcribed genes
H3K9me3	mark of closed heterochromatin
H3K27me3	polycomb associated mark → repressed chromatin

Histone code

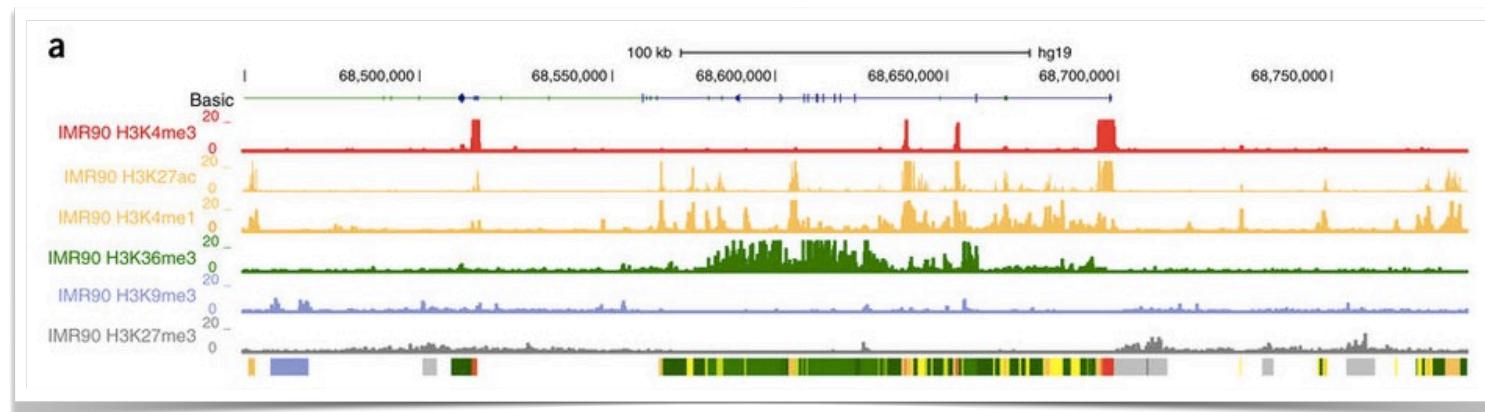


Histone modifications appear to occur in specific combinations related to functional impact → **combinatorial chromatin states**

How can we define/annotate these **chromatin states** ?
→ **Hidden Markov model**

Hidden Markov Model of chromatin states

- There are n (**unobserved**) chromatin states (active promoter, enhancer, repressed chromatin,...)
→ **HIDDEN STATES**
- These states emit specific chromatin features (histone modifications, DNA methylation, accessibility)
→ **OBSERVED DATA**



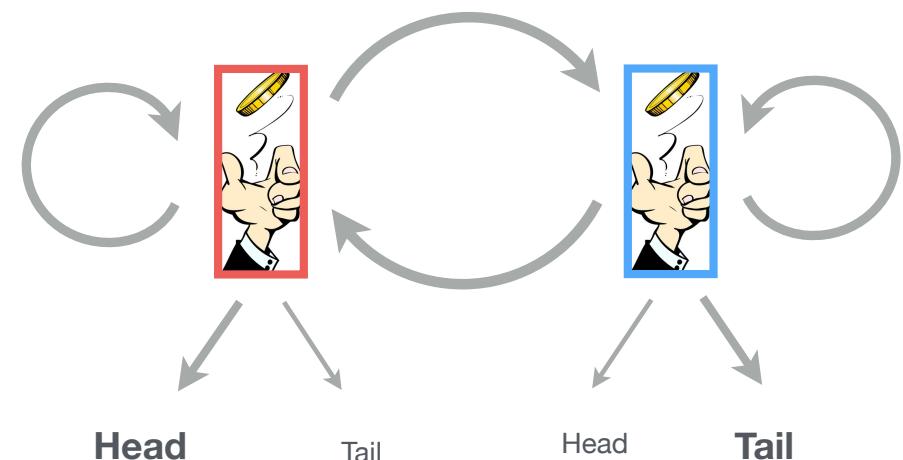
How can we reconstruct the chromatin states from the observed data?

Hidden Markov Models

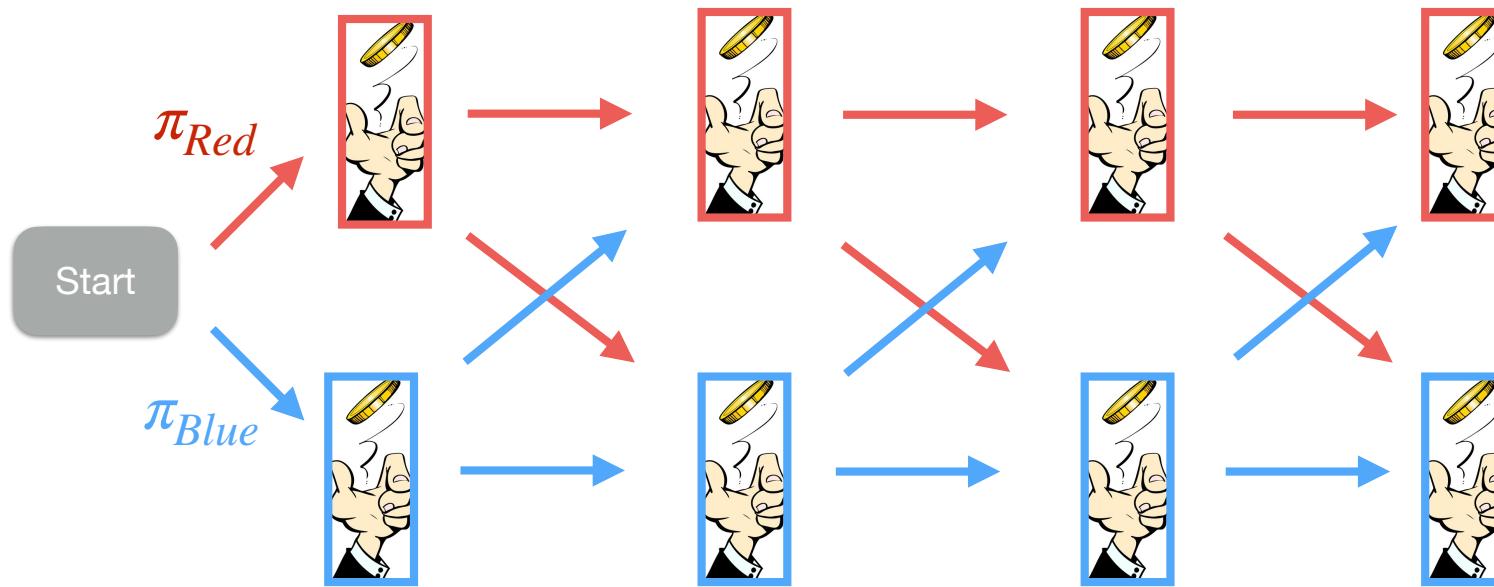
- In Hidden Markov Models, we observe a **sequence of observed events x** , generated by **hidden underlying states z**

$$\begin{array}{cccc} x_1 & x_2 & \dots & x_n \\ \uparrow & \uparrow & \uparrow & \uparrow \\ z_1 & z_2 & \dots & z_n \end{array}$$

- Suppose we have 2 different coins, which are both biased in a different way
- at each step, we
 - choose one of the coins
 - throw it and record the obtained number
 - choose if we
 - ▶ keep the same coin
 - ▶ take the other coin
 - start again ...



Hidden Markov Models



initial probabilities

$$\pi_{Blue} = 0.5$$

$$\pi_{Red} = 0.5$$

B

emission probability

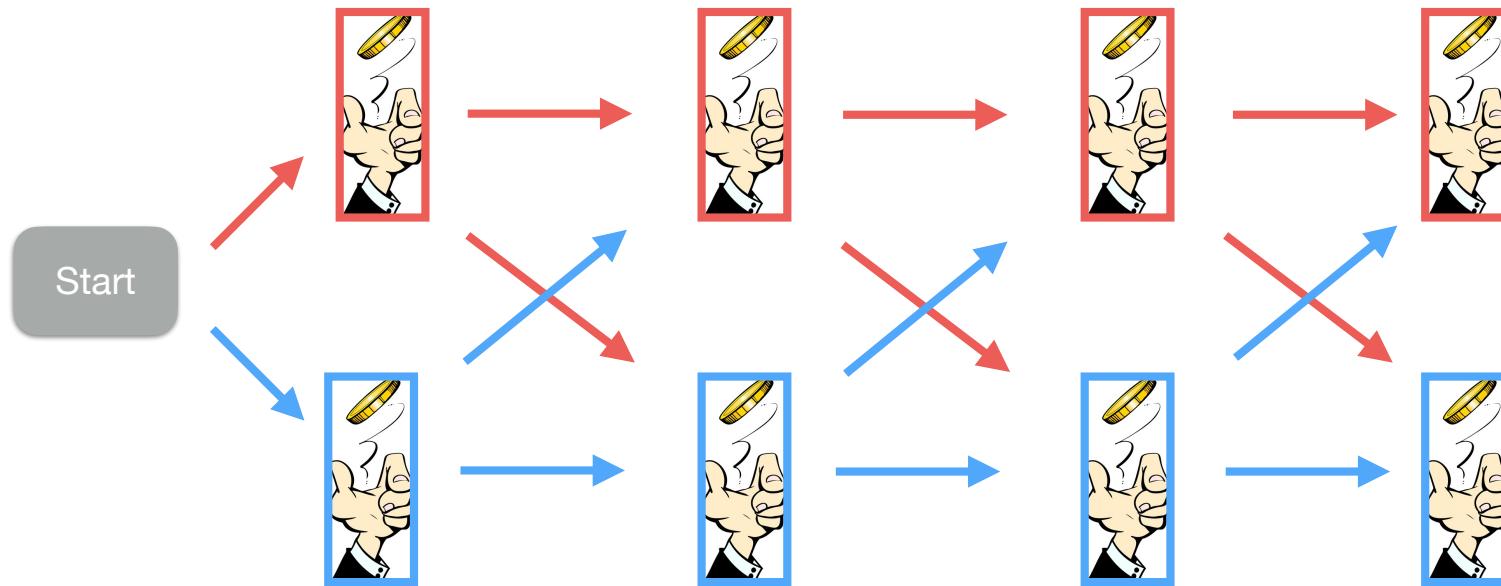
Blue coin : 1 $p=2/3$; 2 $p=1/3$

Red coin: 1 $p=1/3$; 2 $p=2/3$

transition probability

		$p = 0.33$
		$p = 0.67$
		$p = 0.45$
		$p = 0.55$

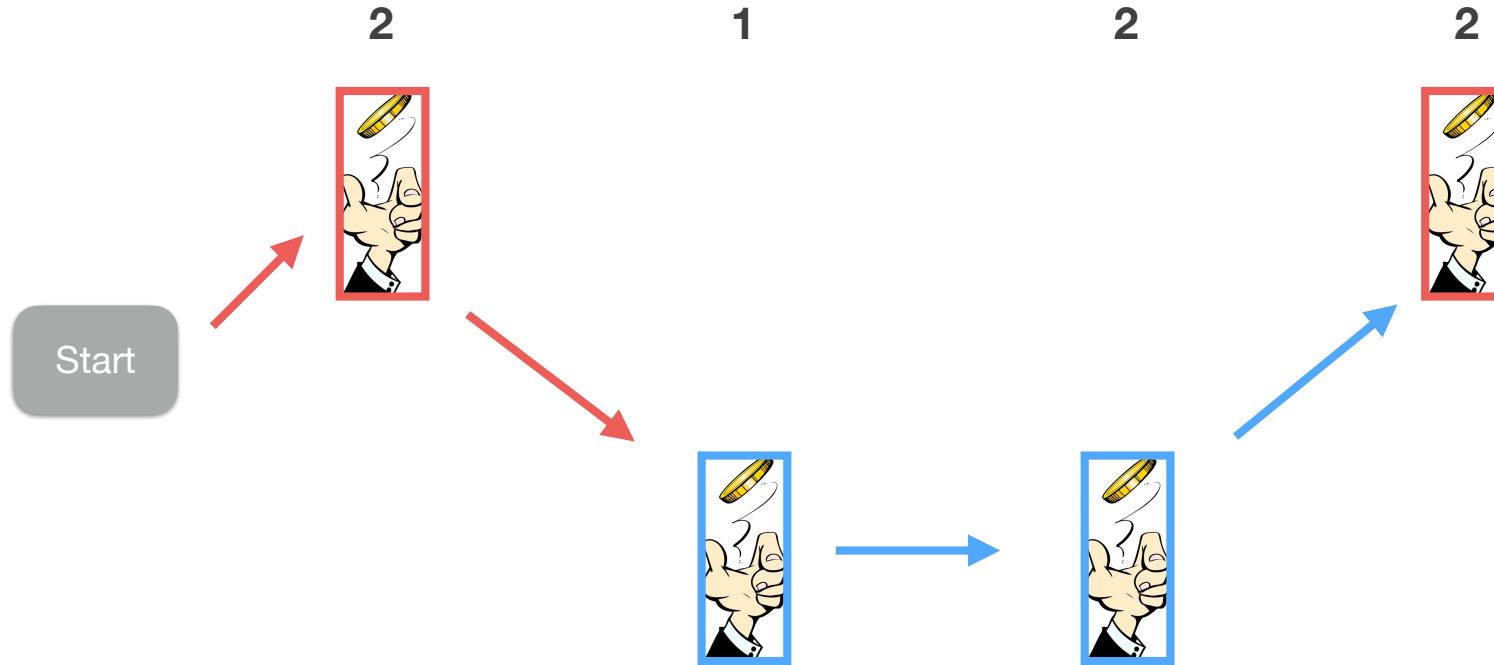
Hidden Markov Models



Markov properties :

- 1. parameters (transition A and emission B)
do not change during the process !!**
- 2. Each throw is independent of the previous one**
- 3. Each transition is independent of the previous one**

Hidden Markov Models



initial probabilities

$$\pi_{Blue} = 0.5$$

$$\pi_{Red} = 0.5$$

B

emission probability

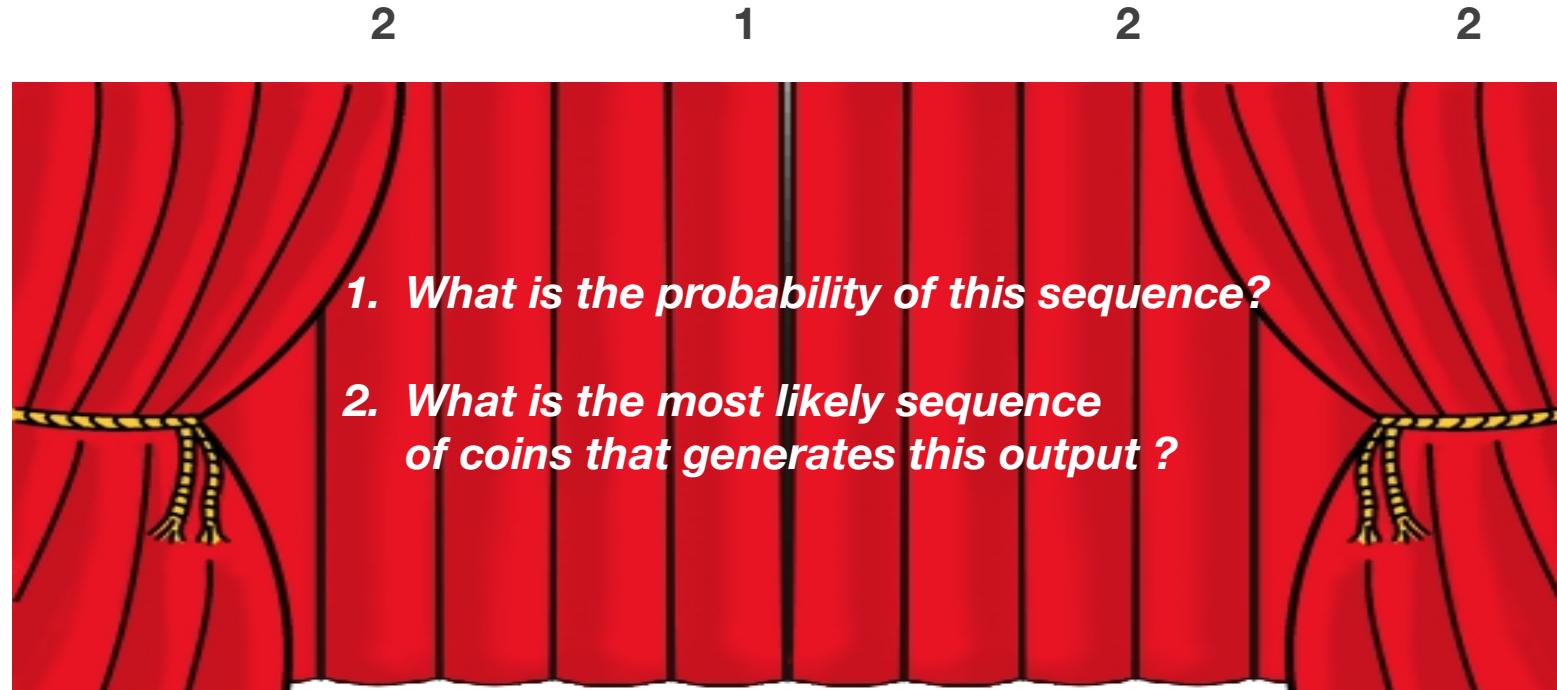
Blue coin : 1 p=2/3 ; 2 p=1/3

Red coin: 1 p=1/3 ; 2 p=2/3

transition probability

 → 	$p = 0.33$
 → 	$p = 0.67$
 → 	$p = 0.45$
 → 	$p = 0.55$

Hidden Markov Models



initial probabilities

$$\pi_{Blue} = 0.5$$

$$\pi_{Red} = 0.5$$

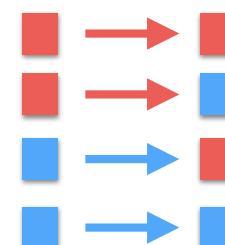
B

emission probability

Blue coin : 1 p=2/3 ; 2 p=1/3

Red coin: 1 p=1/3 ; 2 p=2/3

transition probability



A

Hidden Markov Models

- **Probability of an observed sequence**

- **observed** sequence (length T)
- **unobserved** sequence of states (number of states S):
- Probability of sequence x :

 \vec{x} \vec{z}

$$\begin{aligned}P(\vec{x}; A; B) &= \sum_{\vec{z}} P(\vec{x}; \vec{z}; A; B) \\&= \sum_{\vec{z}} P(\vec{x}|\vec{z}; A; B)P(\vec{z}; A; B)\end{aligned}$$

- Markov assumption:

parameters

A and B remain constant

$$\begin{aligned}P(\vec{x}; A; B) &= \sum_{\vec{z}} P(\vec{x}|\vec{z}; A; B)P(\vec{z}; A; B) \\&= \sum_{\vec{z}} \left(\prod_{t=1}^T P(x_t|z_t; B) \right) \left(\prod_{t=1}^T P(z_t|z_{t-1}; A) \right) \\&= \sum_{\vec{z}} \left(\prod_{t=1}^T B_{z_t, x_t} \right) \left(\prod_{t=1}^T A_{z_{t-1}, z_t} \right)\end{aligned}$$

**emission
probabilities**

**transition
probabilities**

Hidden Markov Models

$$P(\vec{z} | A; B) = \text{Red} \rightarrow \text{Red} \rightarrow \text{Red} \rightarrow \text{Red} = \frac{2}{3} * 0.33 * \frac{1}{3} * \frac{2}{3} * 0.33 * \frac{2}{3} * 0.33 * \frac{2}{3} = 0.0036$$

$$P(\vec{x} | \vec{z}; A; B) = \text{Blue} \rightarrow \text{Blue} \rightarrow \text{Red} \rightarrow \text{Red} \rightarrow \text{Blue} = \frac{1}{3} * 0.55 * \frac{2}{3} * 0.45 * \frac{2}{3} * 0.67 * \frac{1}{3} = 0.0081$$

⋮
⋮

initial probabilities

$$\pi_{Blue} = 0.5$$

$$\pi_{Red} = 0.5$$

B

emission probability

Blue coin : 1 p=2/3 ; 2 p=1/3

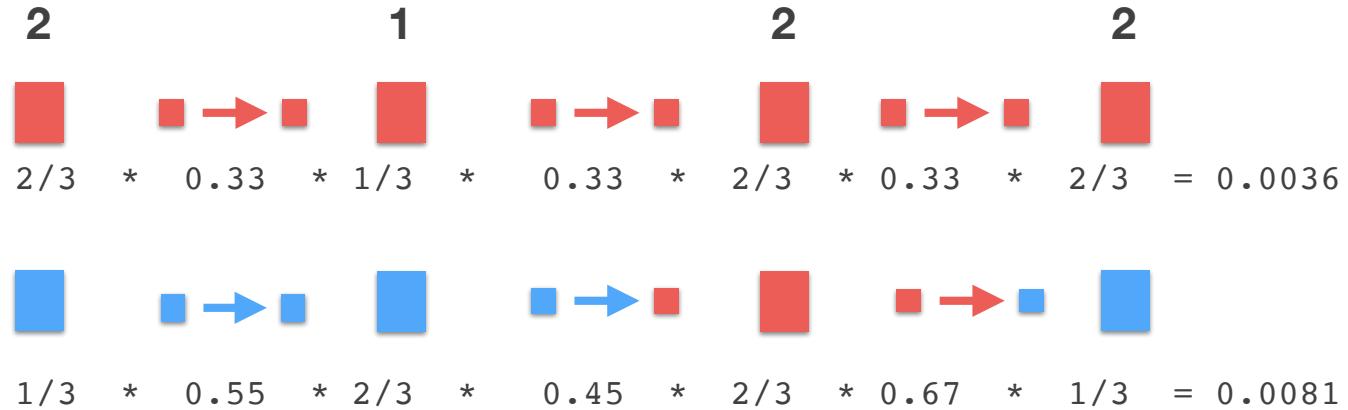
Red coin: 1 p=1/3 ; 2 p=2/3

transition probability

	$p = 0.33$
	$p = 0.67$
	$p = 0.45$
	$p = 0.55$

A

Hidden Markov Models

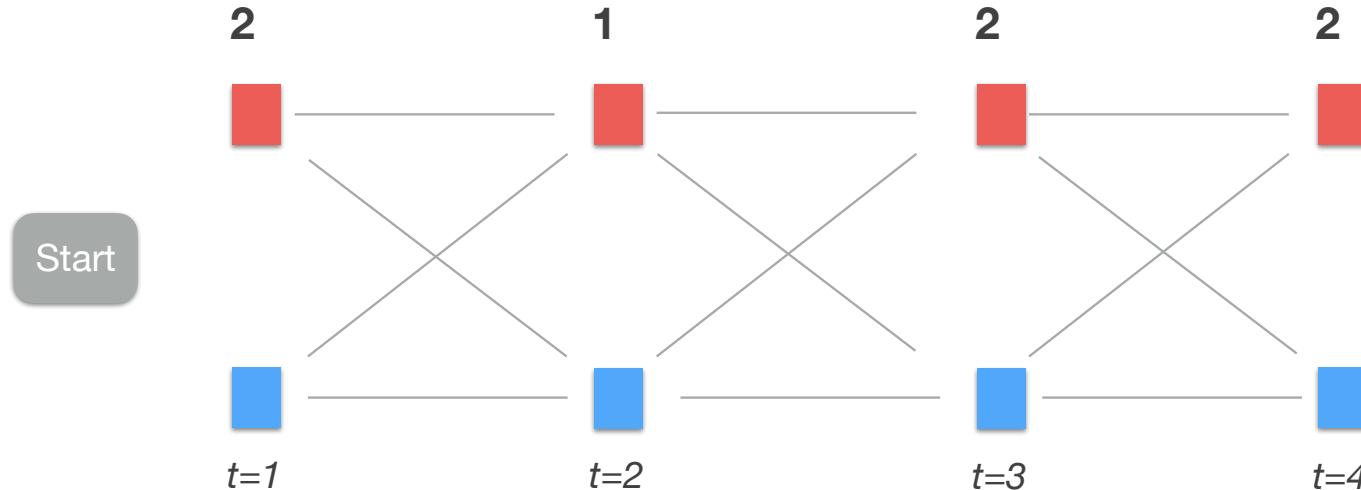


How many possible sequences ? $\rightarrow 2^4 = 16$
In general : $\rightarrow o(|S|^T)$



(Chromatin states : S = 18 ; T = 10^7 ...)

Hidden Markov Models



To compute the probability of an observed sequence, we use the **Forward algorithm** (dynamical programming)

For each time point t_0 and each possible state at $t=t_0$, compute the total probability for the observed sequence until t_0

$P(2-1-2, \square)$ = probability of observing the sequence 2-1-2 and being in the red state at $t=3$

Hidden Markov models

- **Forward algorithm :**

- probability of observing sequences until t , and being in state i at t

$$\alpha_i(t) = P(x_1, x_2, \dots, x_t, z_t = s_i; A; B)$$

- Total probability:

$$P(\vec{x}; A; B) = \sum_{i=1}^{|S|} \alpha_i(T)$$

- Dynamic programming:

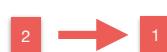
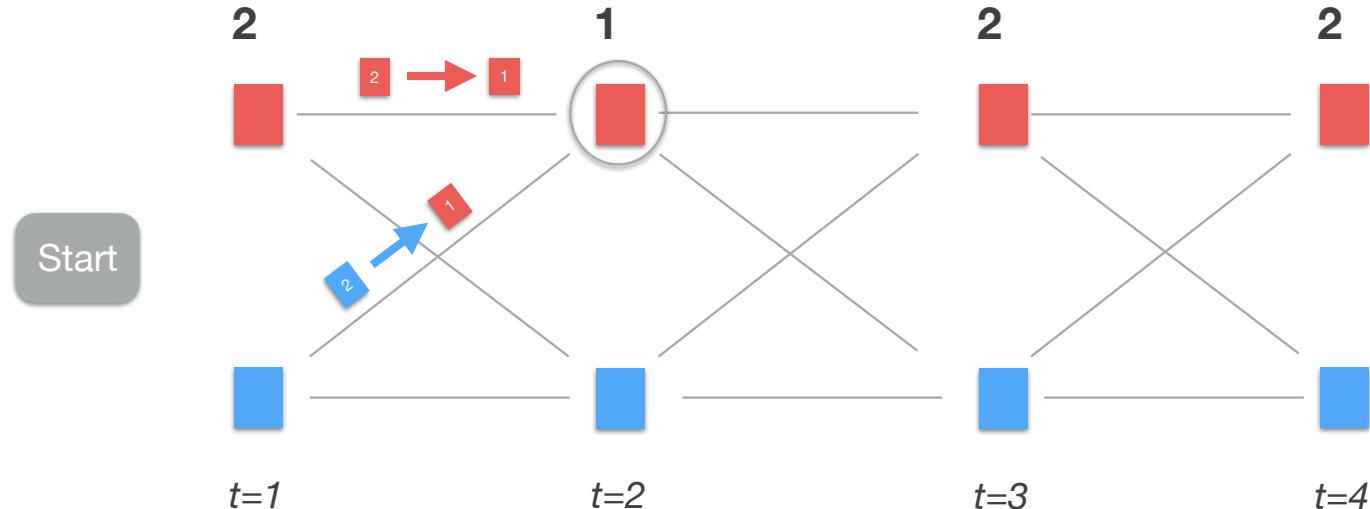
Algorithm 1 Forward Procedure for computing $\alpha_i(t)$

1. Base case: $\alpha_i(0) = A_{0,i}$, $i = 1..|S|$
2. Recursion: $\alpha_j(t) = \sum_{i=1}^{|S|} \alpha_i(t-1) A_{ij} B_{j|x_t}$, $j = 1..|S|$, $t = 1..T$

- Complexity: $O(|S|^2 T)$



Hidden Markov Models



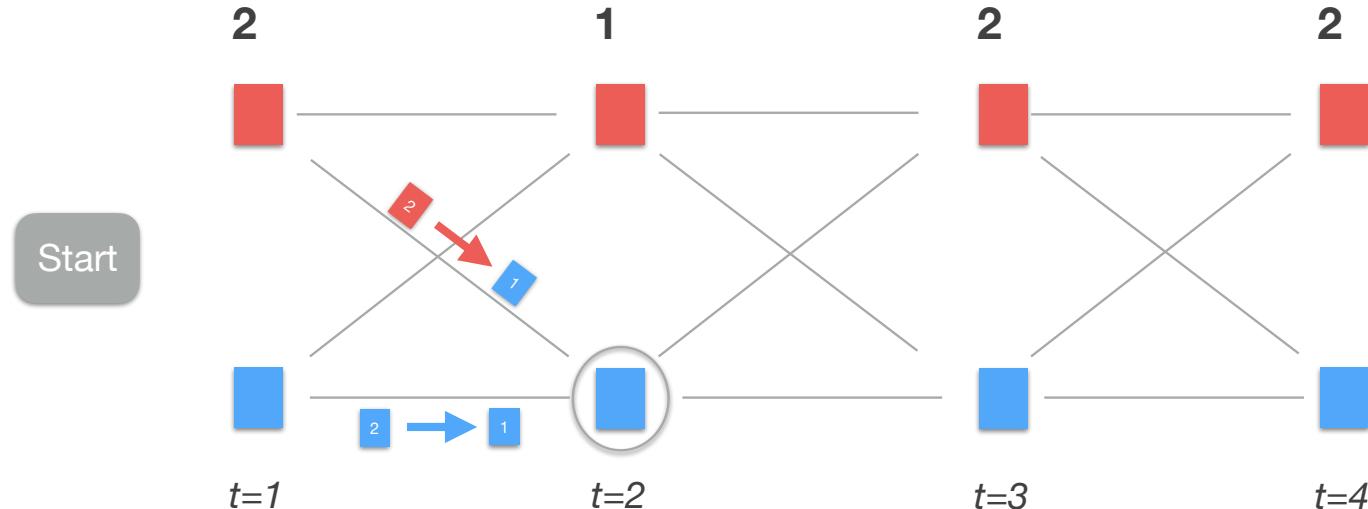
Probability red=2 * probability red \rightarrow red * probability red=1
 $= 2/3 * 1/3 * 1/3 = 0.074$



Probability blue=2 * probability blue \rightarrow red * probability red=1
 $= 1/3 * 0.45 * 1/3 = 0.05$

$$\text{Probability}(2-1, \text{red}) = 0.074 + 0.05 = 0.124$$

Hidden Markov Models



 → 

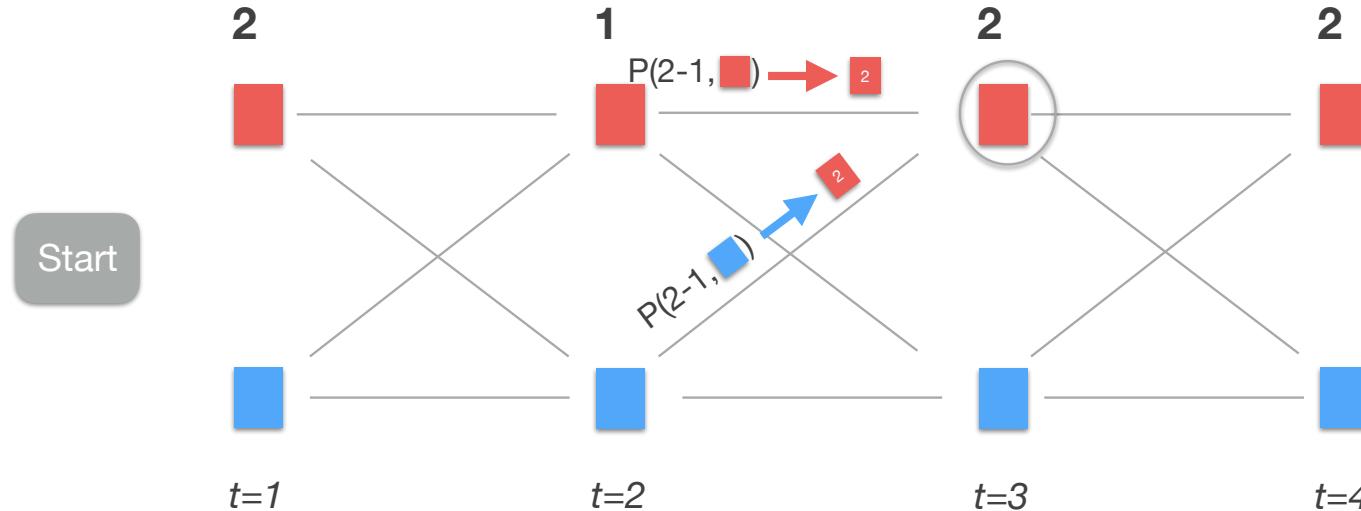
Probability blue=2 * probability blue→ blue * probability blue=1
 $= 1/3 * 0.55 * 2/3 = 0.12$

 → 

Probability red=2 * probability red→ blue * probability blue=1
 $= 2/3 * 2/3 * 2/3 = 0.296$

Probability(2-1, ) = $0.12 + 0.296 = 0.416$

Hidden Markov Models



$$P(2-1, \text{red}) \rightarrow \text{red}_2$$

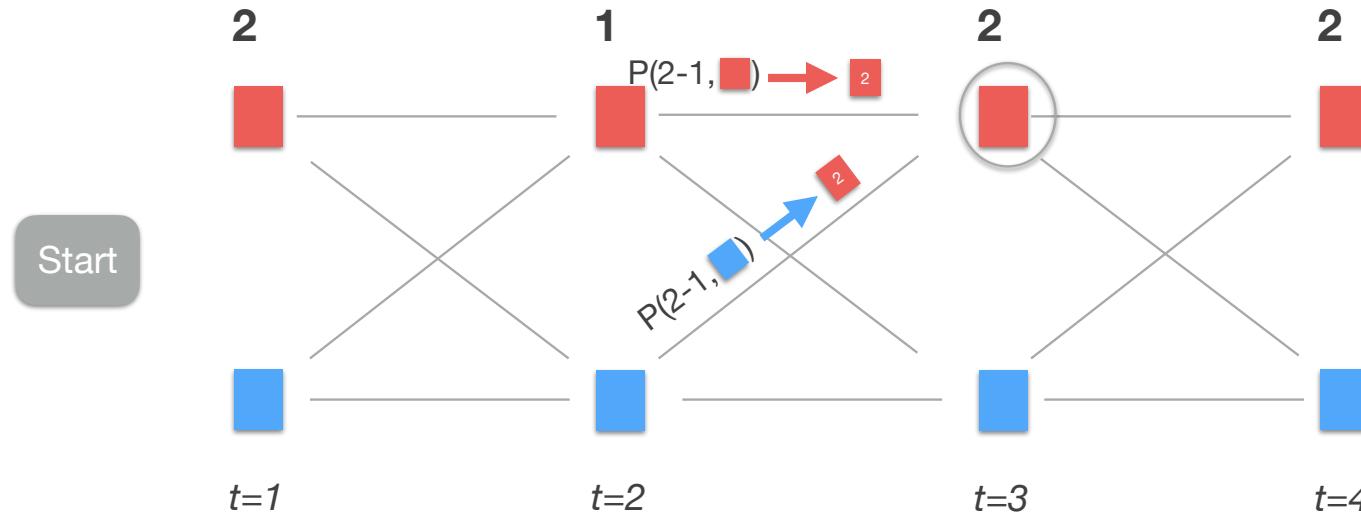
$$\begin{aligned} & P(2-1, \text{red}) * \text{probability red} \rightarrow \text{red} * \text{probability red}=2 \\ & = 0.124 * 1/3 * 2/3 = 0.0276 \end{aligned}$$

$$P(2-1, \text{blue}) \rightarrow \text{blue}_2$$

$$\begin{aligned} & P(2-1, \text{blue}) * \text{probability blue} \rightarrow \text{red} * \text{probability red}=2 \\ & = 0.416 * 0.45 * 2/3 = 0.1249 \end{aligned}$$

$$\text{Probability}(2-1-2, \text{red}) = 0.0276 + 0.1249 = 0.1525$$

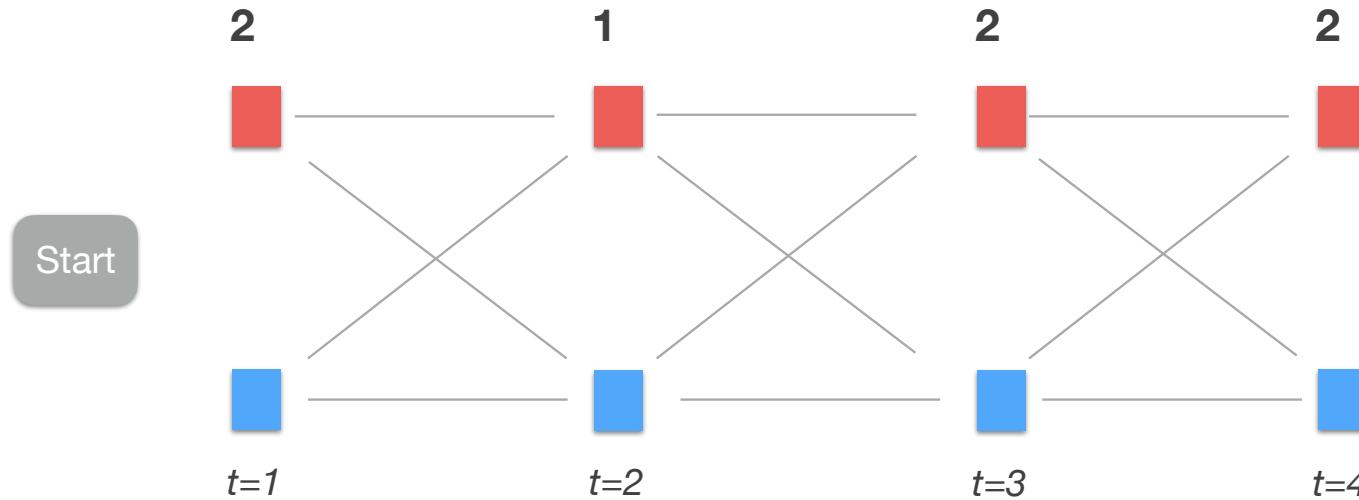
Hidden Markov Models



- at each time point, we perform S^*S operations (S = number of possible states)
- all together: S^*S^*T operations → complexity

$$o(S^2 \cdot T)$$

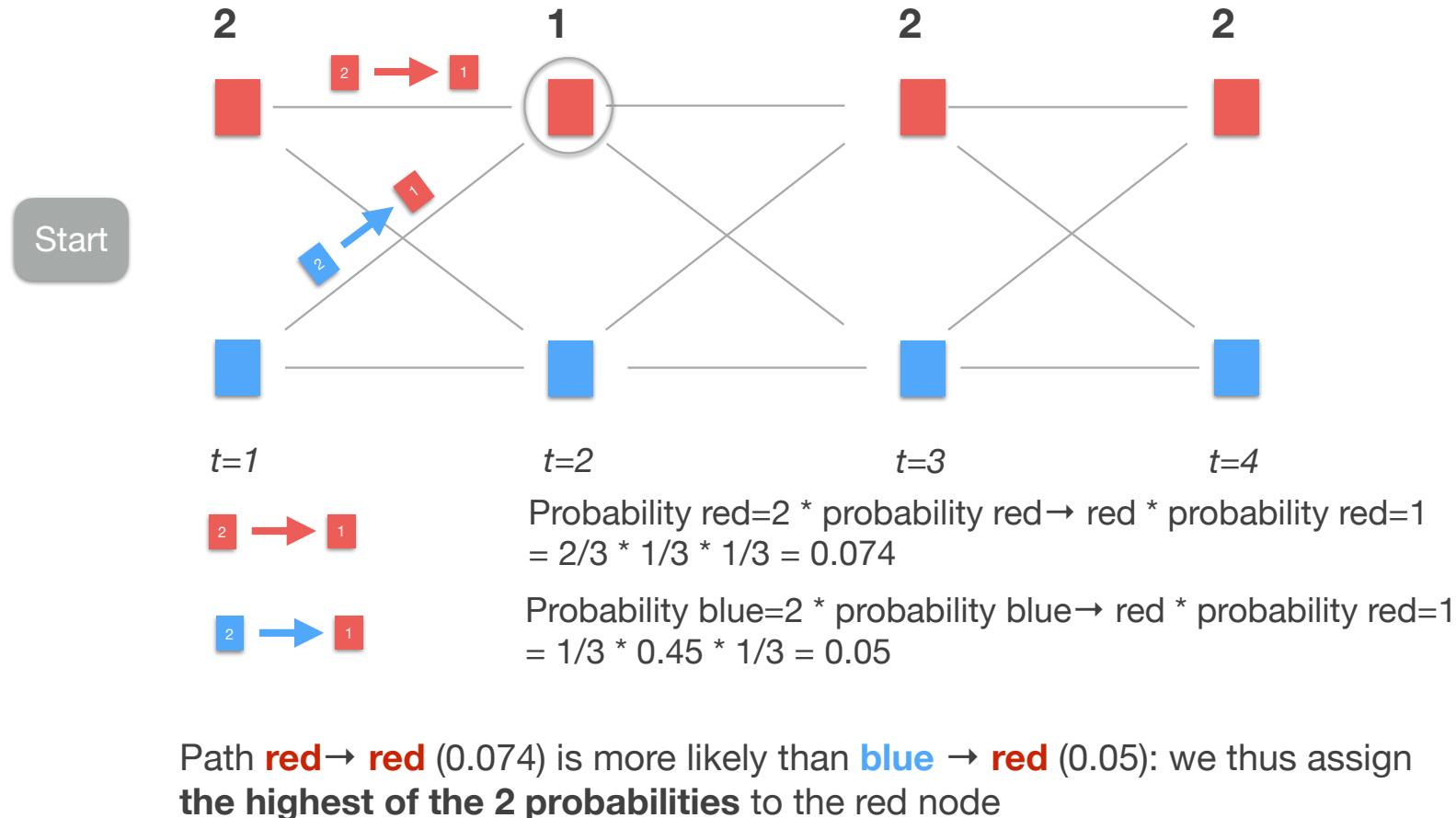
Viterbi algorithm



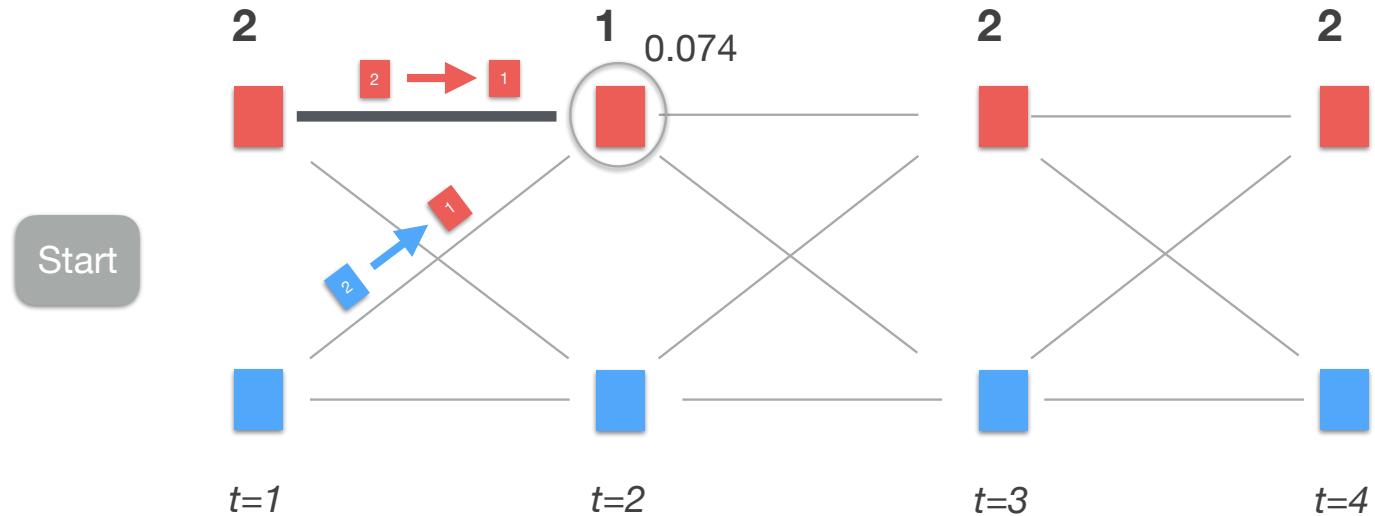
What is the most likely sequence of states leading to the observed sequence ? → **Viterbi algorithm**

For each time point and each possible state, **record the most probable path leading to this state** at this time point

Viterbi algorithm



Viterbi algorithm



Probability red=2 * probability red→ red * probability red=1
 $= 2/3 * 1/3 * 1/3 = 0.074$

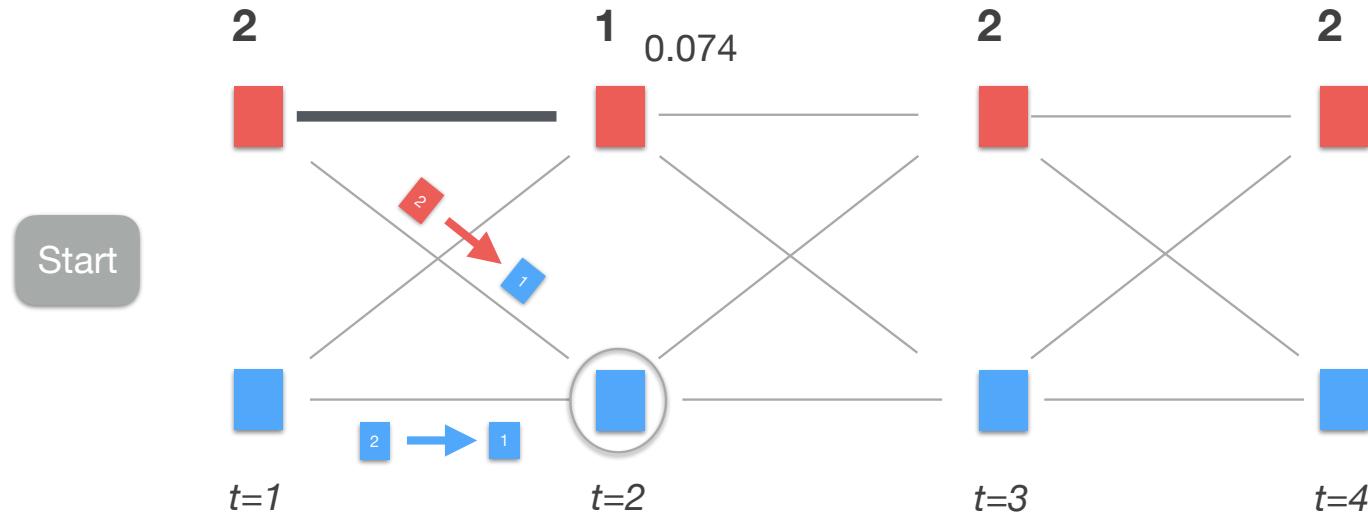


Probability blue=2 * probability blue→ red * probability red=1
 $= 1/3 * 0.45 * 1/3 = 0.05$



This is unlike the Forward algorithm, where we would assign the sum of the 2 probabilities ($0.074+0.05 = 0.124$) to the red node!! (compare with slide 129)

Viterbi algorithm

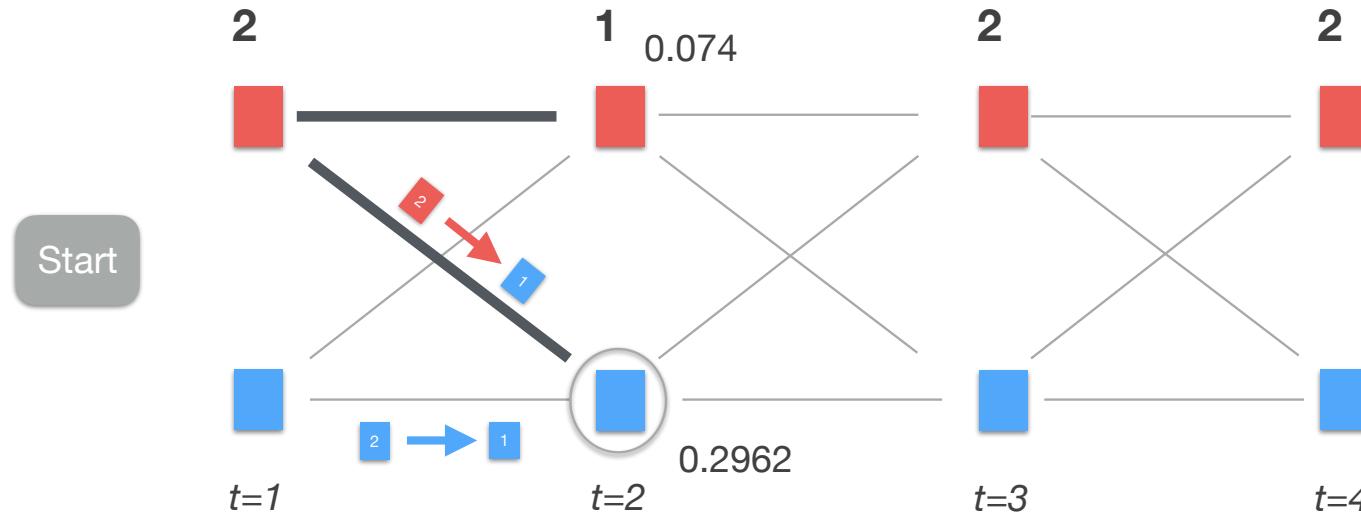


Probability red=2 * probability red \rightarrow blue * probability blue=1
 $= 2/3 * 2/3 * 2/3 = 0.2962$



Probability blue=2 * probability blue \rightarrow blue * probability blue=1
 $= 1/3 * 0.55 * 2/3 = 0.1222$

Viterbi algorithm

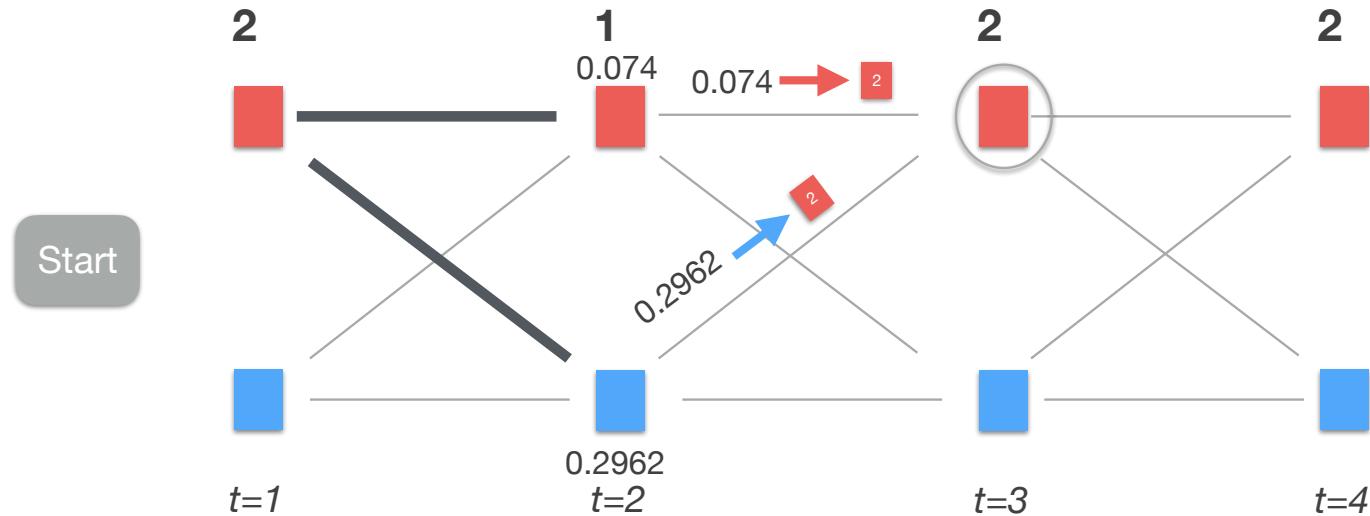


Probability red=2 * probability red→ blue * probability blue=1
 $= 2/3 * 2/3 * 2/3 = 0.2962$



Probability blue=2 * probability blue→ blue * probability blue=1
 $= 1/3 * 0.55 * 2/3 = 0.1222$

Viterbi algorithm



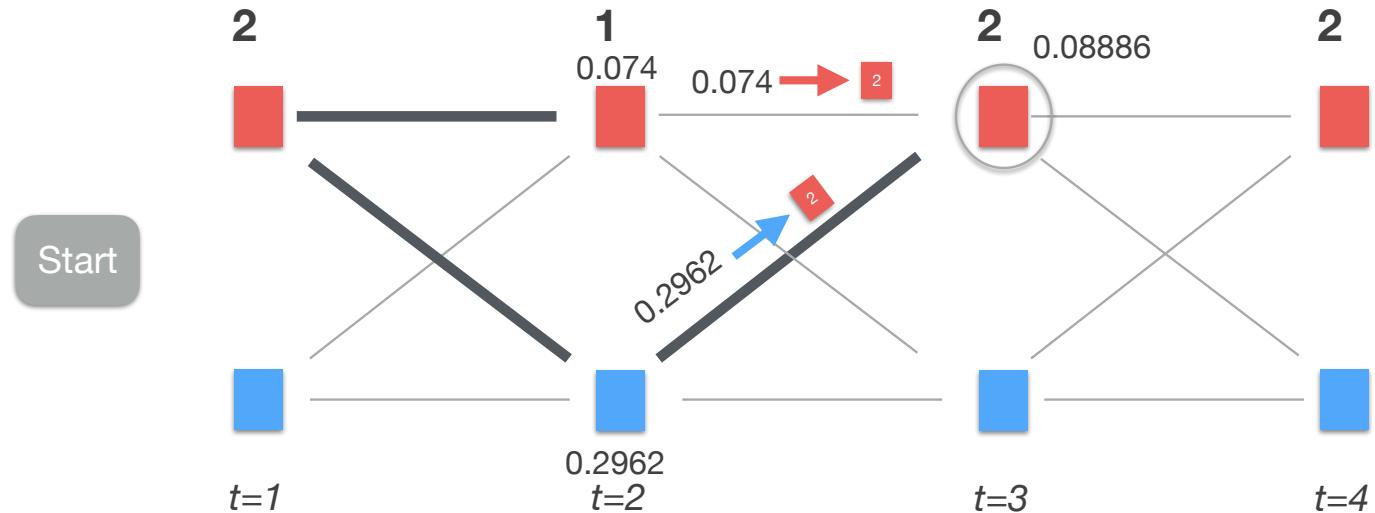
0.074 → 2

$$0.074 * \text{probability red} \rightarrow \text{red} * \text{probability red}=2 \\ = 0.074 * 1/3 * 2/3 = 0.01644$$

0.2962 → 2

$$0.2962 * \text{probability blue} \rightarrow \text{red} * \text{probability red}=2 \\ = 0.2962 * 0.45 * 2/3 = 0.08886$$

Viterbi algorithm



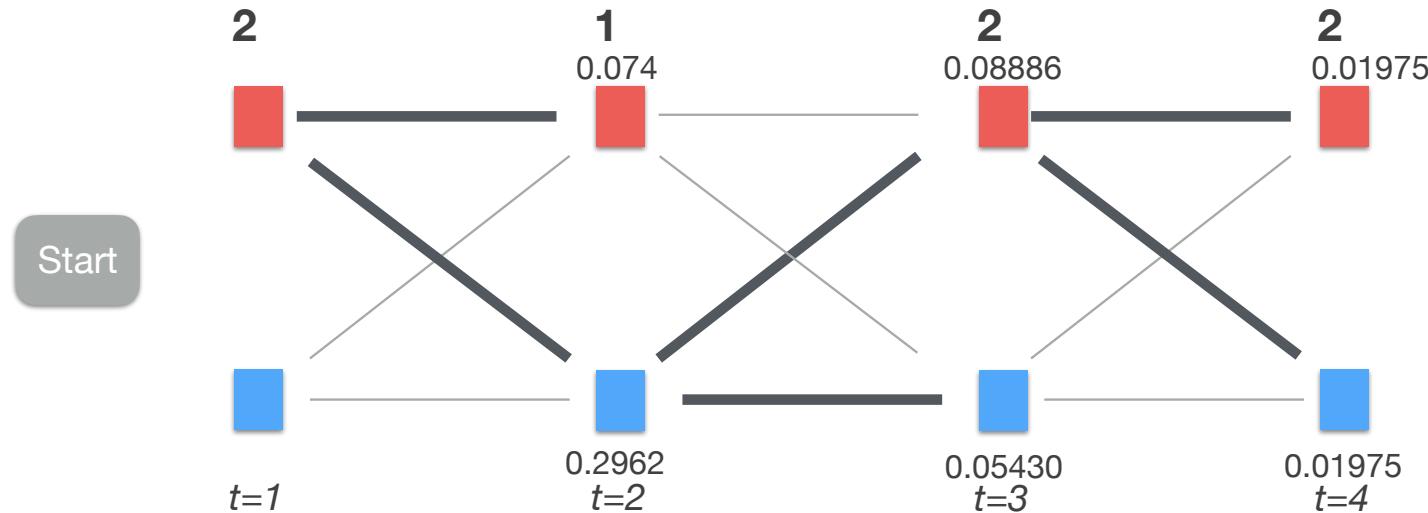
0.074 → 2

$$0.074 * \text{probability red} \rightarrow \text{red} * \text{probability red}=2 \\ = 0.074 * 1/3 * 2/3 = 0.01644$$

0.2962 → 2

$$0.2962 * \text{probability blue} \rightarrow \text{red} * \text{probability red}=2 \\ = 0.2962 * 0.45 * 2/3 = 0.08886$$

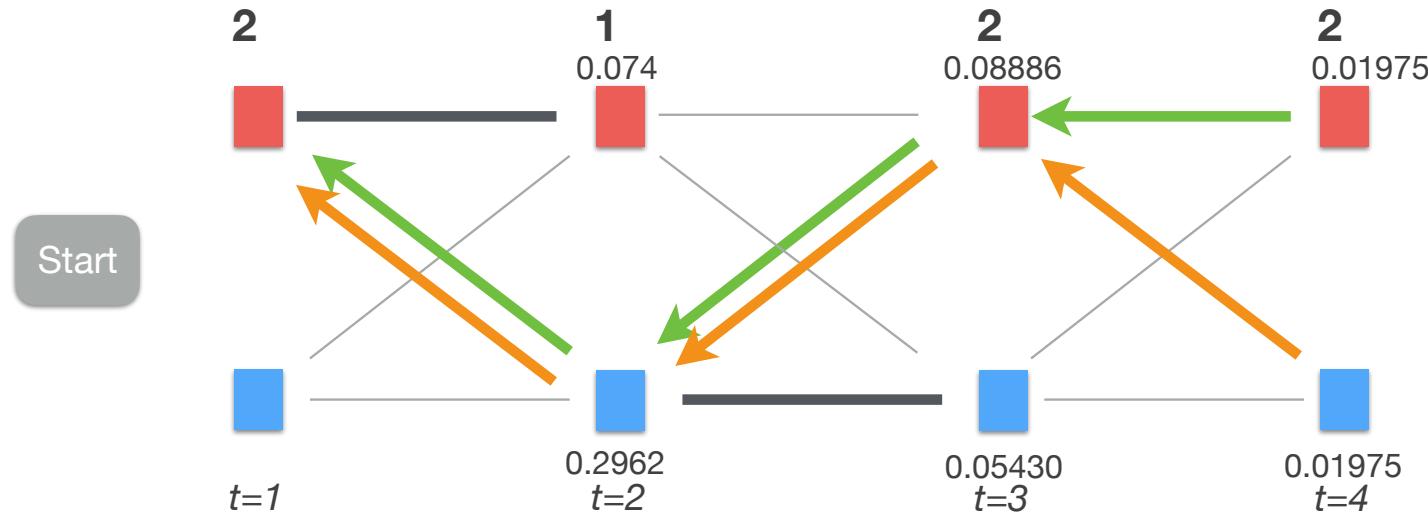
Viterbi algorithm



We have two end nodes with equal probability!
Hence, we have 2 optimal path obtained by backtracking
from the 2 possible optimal final states!

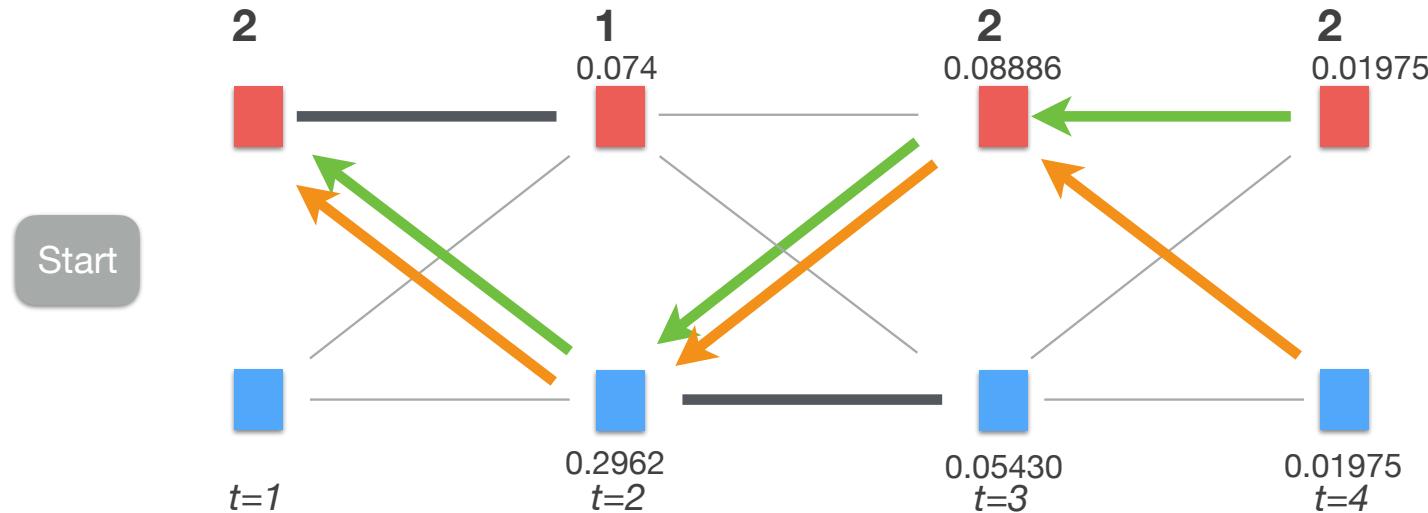
Generally, we would have only one end node with highest probability
hence a single backtracking path!

Viterbi algorithm



We have **two final states with equal probability!**
Hence, we have 2 optimal path (**green** and **orange**)
obtained by backtracking from the 2 possible optimal final states!

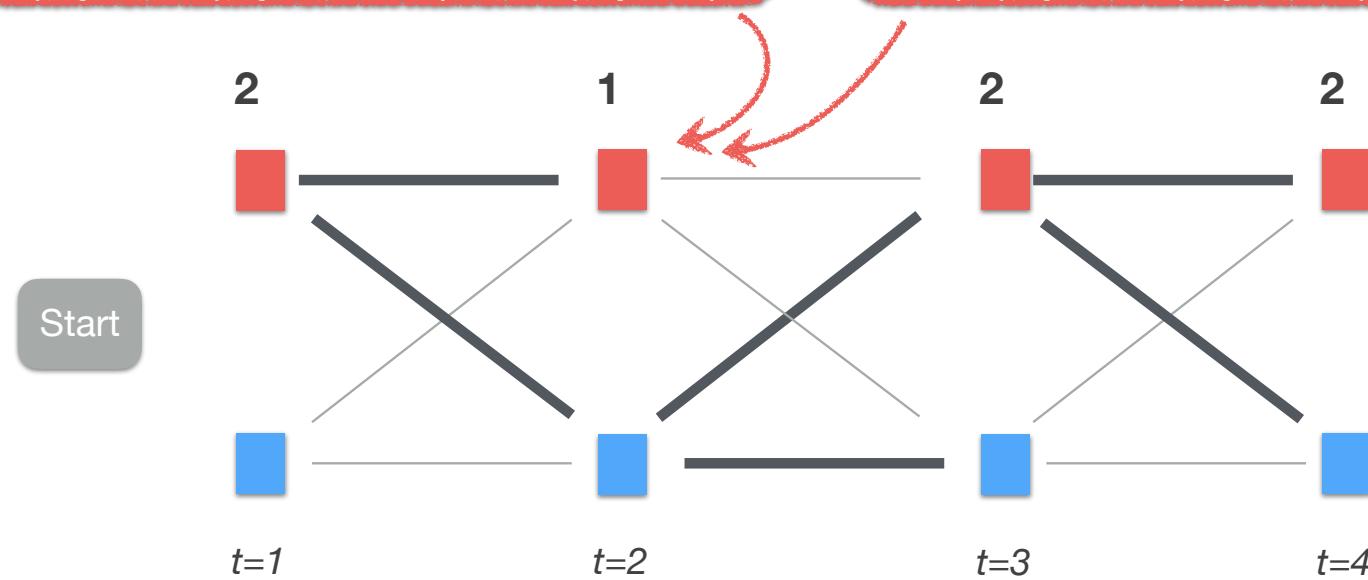
Viterbi algorithm



Summary: difference Forward / Viterbi

Forward: assign to the node at t the sum of the probabilities of all path leading to that node from $t-1$

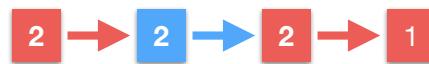
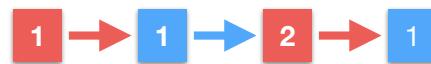
Viterbi: assign to the node at t the probability of the most likely path leading to that node from $t-1$ + backtracking



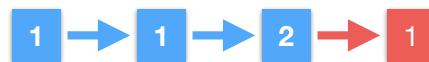
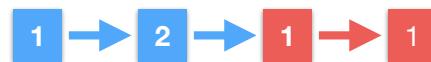
HMM - learning parameters



- So far, we assumed that the parameters A and B are given...
How can we determine them ?
→ use a set of **training sequences** with **known states**



etc ...

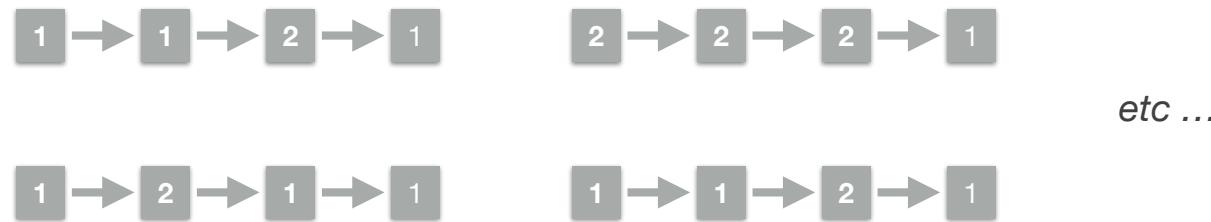


$$P(\boxed{1}) = \frac{\boxed{1}}{\boxed{1} + \boxed{2}} = 5/8$$

$$P(\boxed{\quad} \rightarrow \boxed{\quad}) = \frac{\boxed{\quad} \rightarrow \boxed{\quad}}{\boxed{\quad} \rightarrow \boxed{\quad} + \boxed{\quad} \rightarrow \boxed{\quad}} = 2/5 \quad \text{etc ...}$$

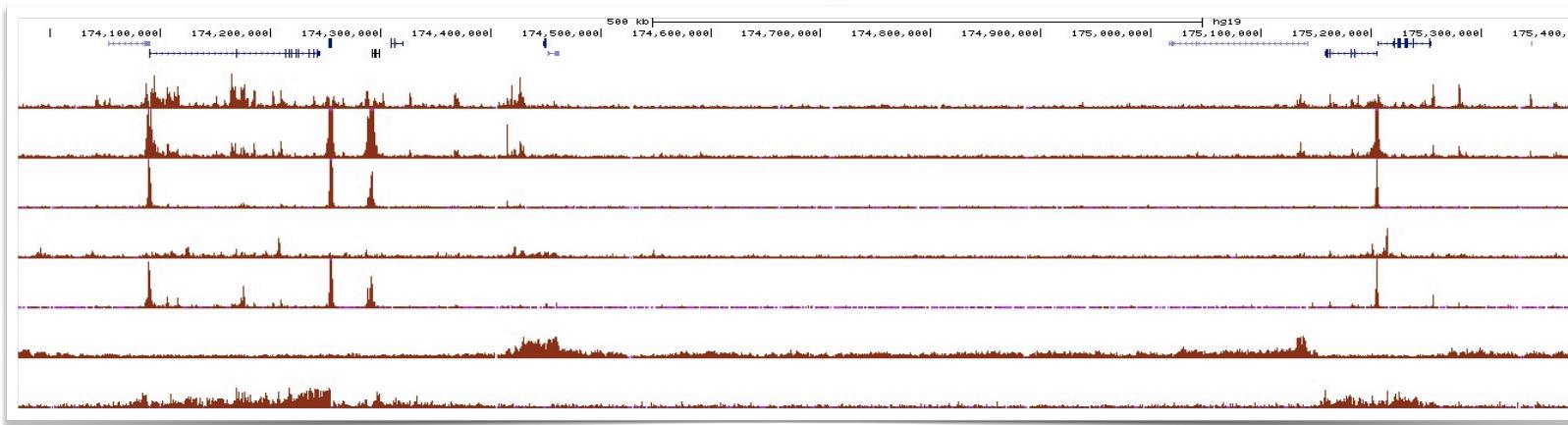
HMM - learning parameters

- So far, we assumed that the parameters A and B are given...
How can we determine them ?
→ use a set of **training sequences** with **unknown states**



- Initialize randomly parameters; update these parameters iteratively until convergence
(*Baum-Welch* algorithm, based on **Expectation Maximization**)

HMM and Chromatin states



n histone marks;

Binarize the data in windows of 200bp (using appropriate thresholds)

[ChromHMM]

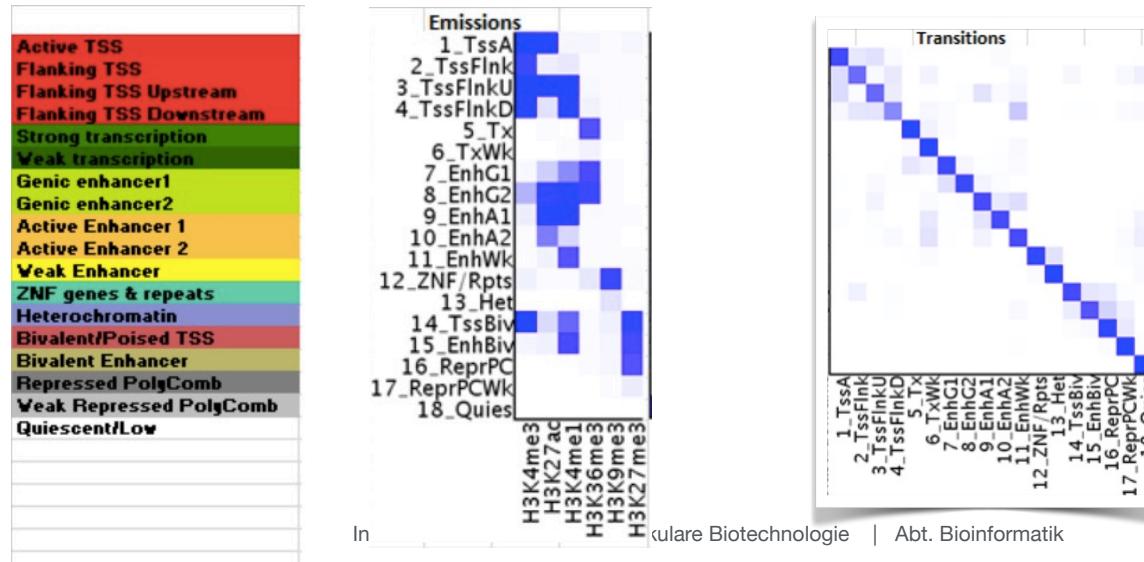
Each observation is a binary vector of length n in a genomic window of 200bp

Chromatin states

We postulate the existence of k (= 15 or 18) chromatin states
Each state has different **emission probability** for each histone mark ...



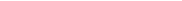
Emission and transition probabilities are learned from training data:



Chromatin states

- Given that we observe

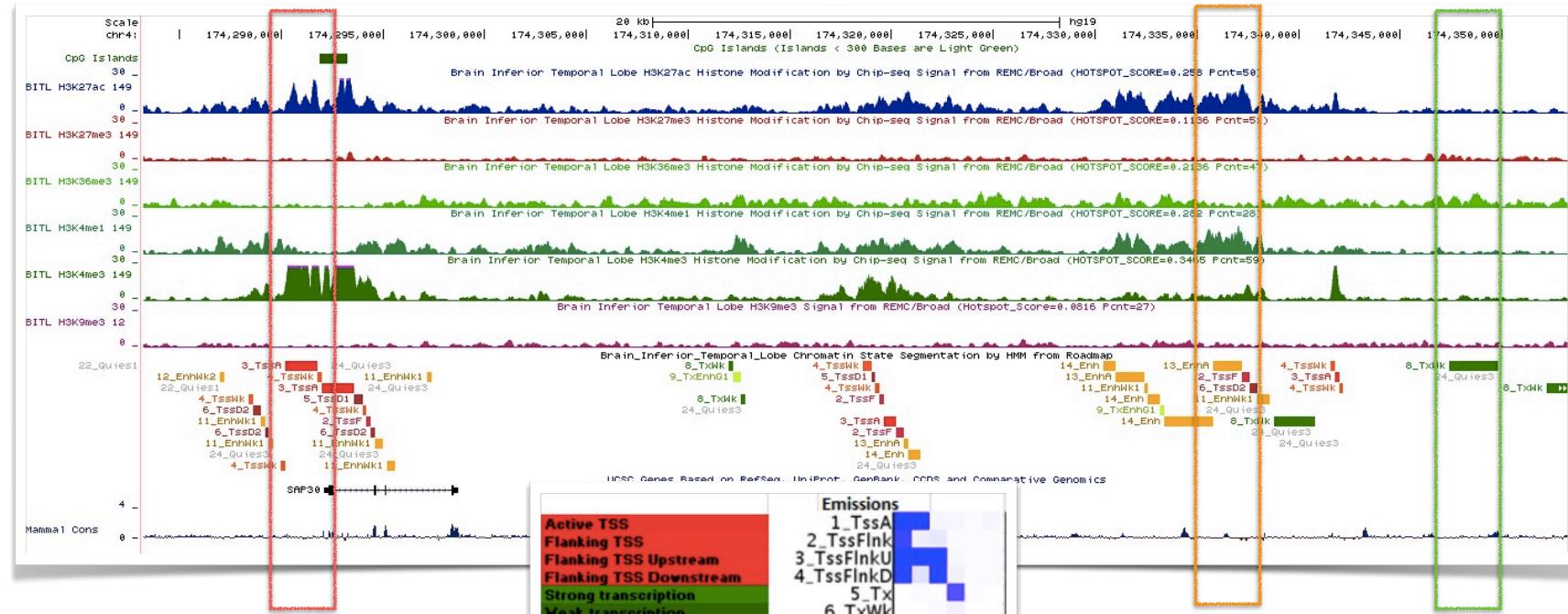
... what is the most likely sequence of chromatin states leading to these observed values ? → ***Viterbi algorithm***

states : 

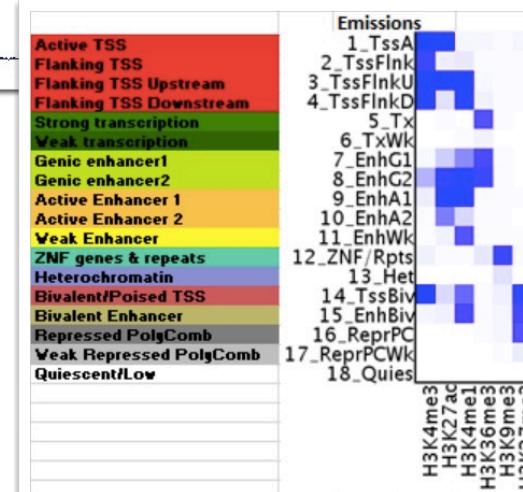
purple box

Page 10

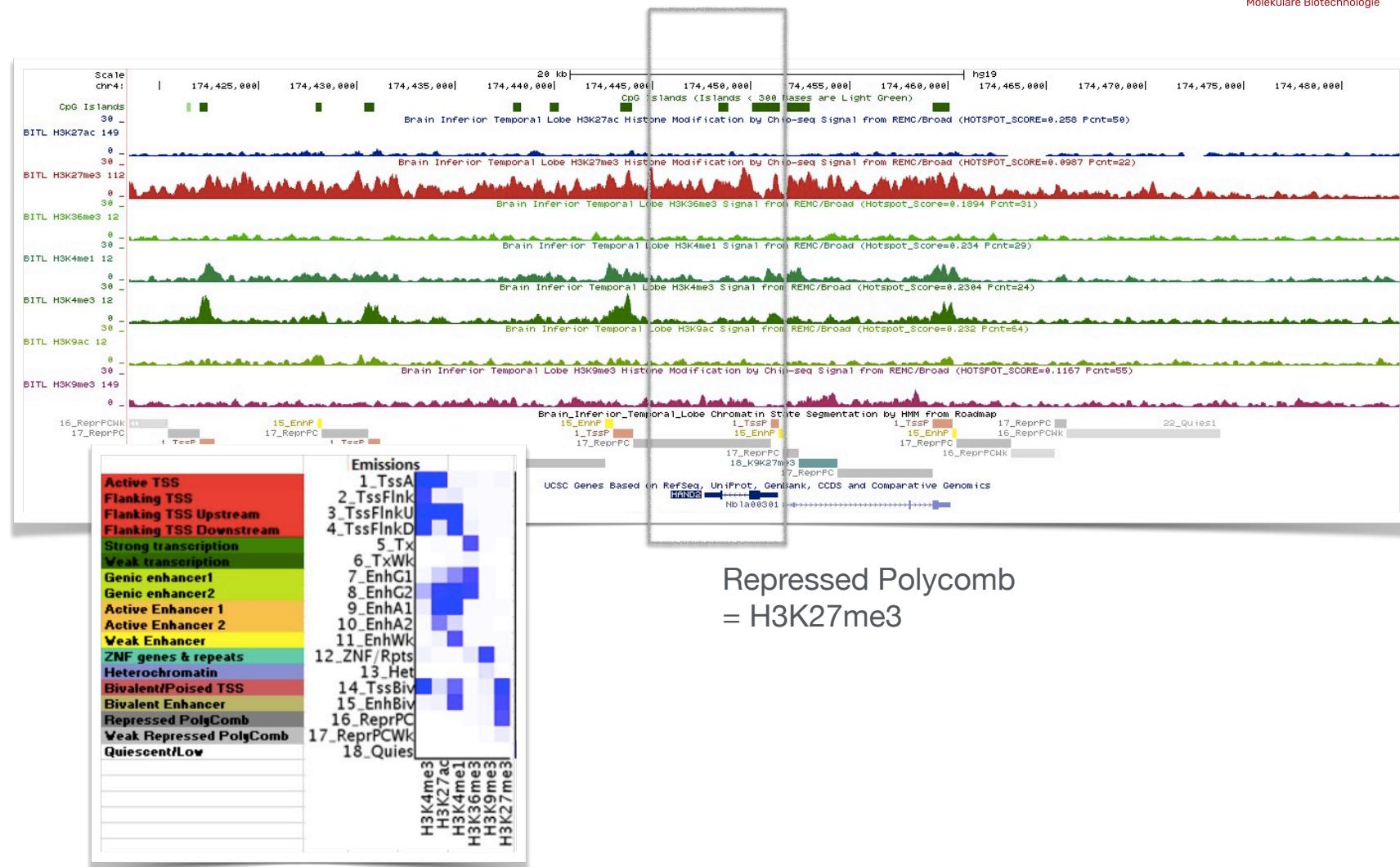
Chromatin states



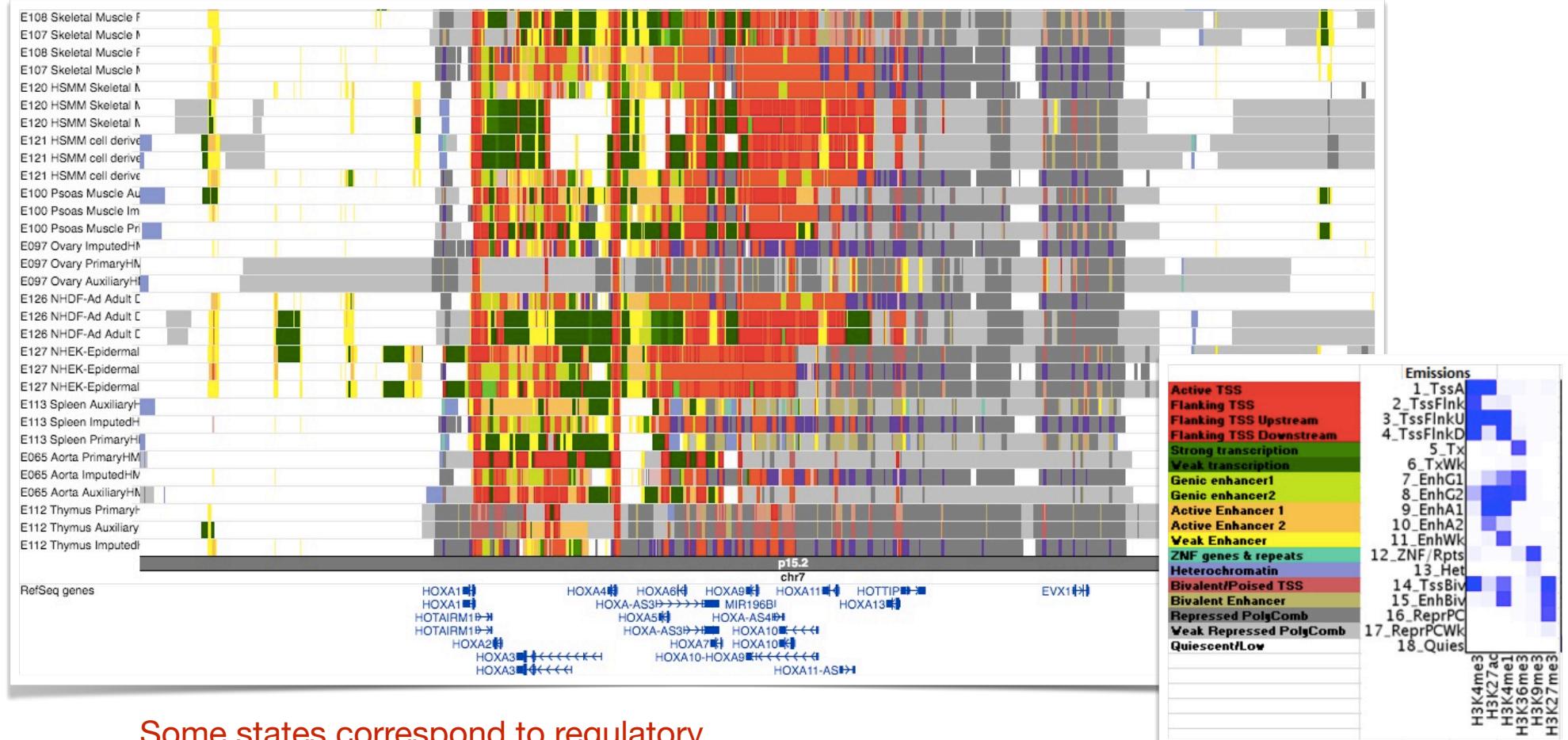
Active Transcription
Start Site
= H3K4me3 + H3K27ac



Chromatin states



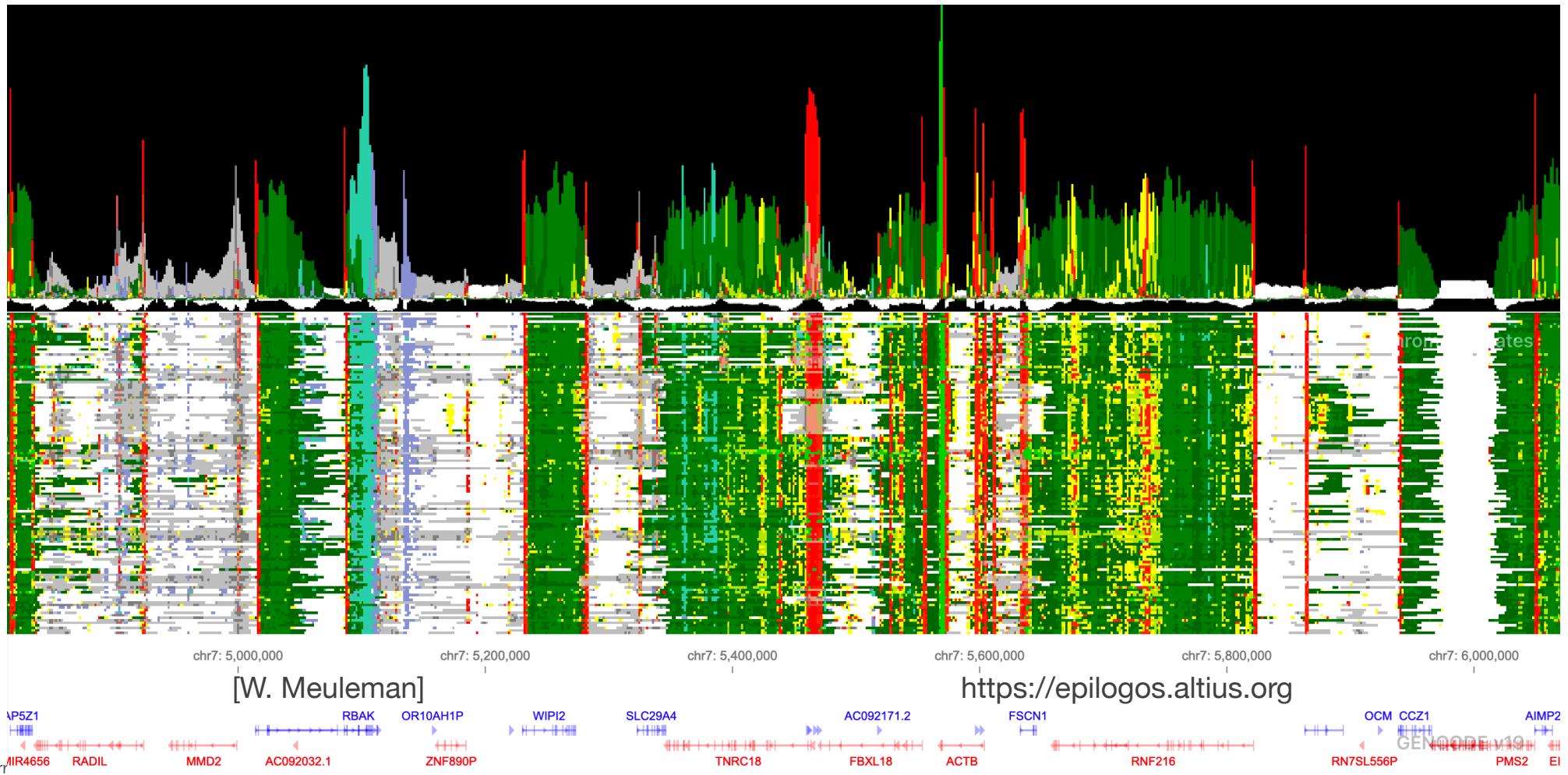
Roadmap chromatin segmentation in different human adult tissues



Some states correspond to regulatory regions (active and poised enhancers)
→ motif search can be restricted to these regions

<http://epigenomegateway.wustl.edu>

EpiLogos



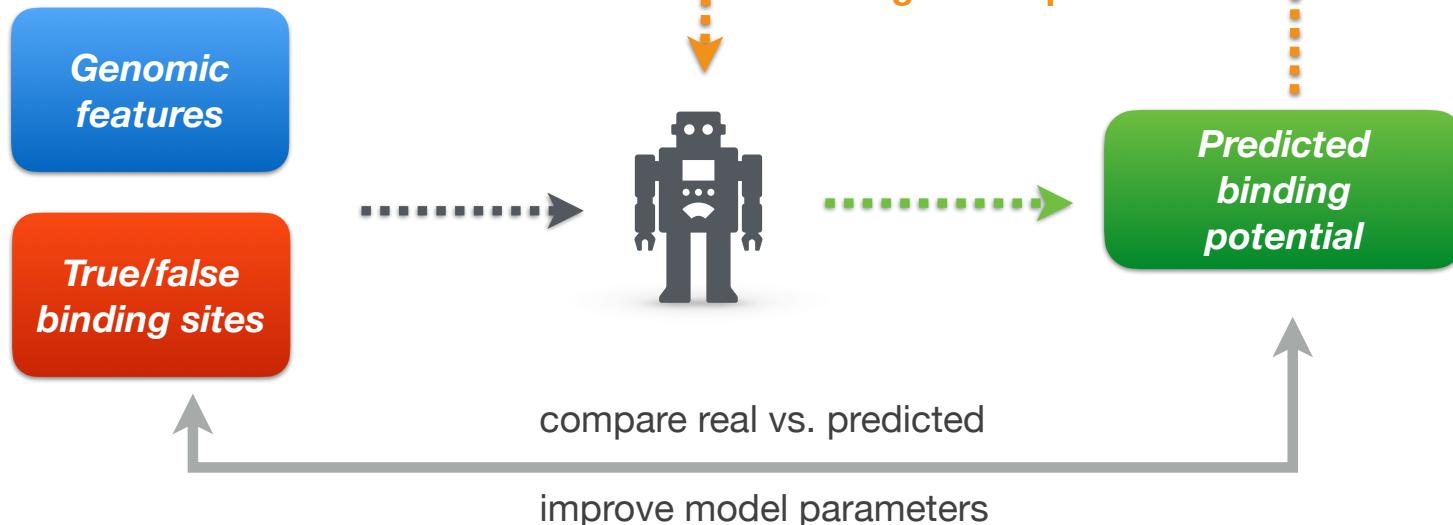
Machine-learning approach to TFBS prediction

- Define a **General Binding Potential** (GBP) for each base in the genome using a **Logistic Regression Approach** (LRA) based on a set of features
→ *What features discriminate best bound/unbound bases ?*
- Training set
 - **positive** : center of ChIP peaks for 14 TF datasets in various cell lines
 - **negative** : randomly sampled bases (50x more negative sequences)
- Features
 - sequence features : evolutionary conservation
 - distance to TSS
 - density of ChIP-seq reads for various histone marks and Pol2
 - presence of DNase I hypersensitive region
 - ...

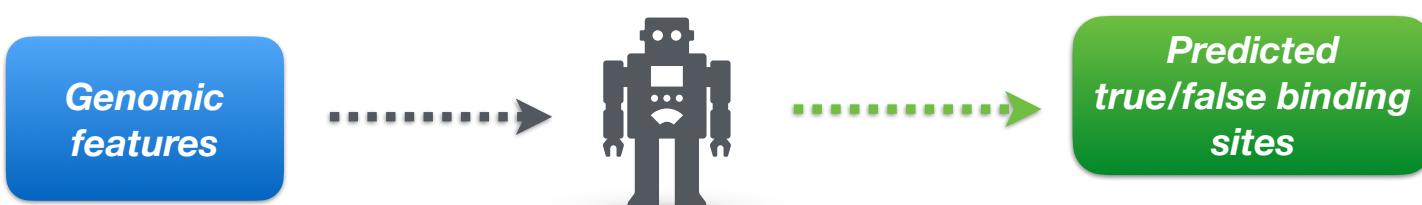
29 features
no PWM data

Principle of the model

1. Training



2. Prediction



Machine-learning approach to TFBS prediction

evolutionary
features

sequence
properties

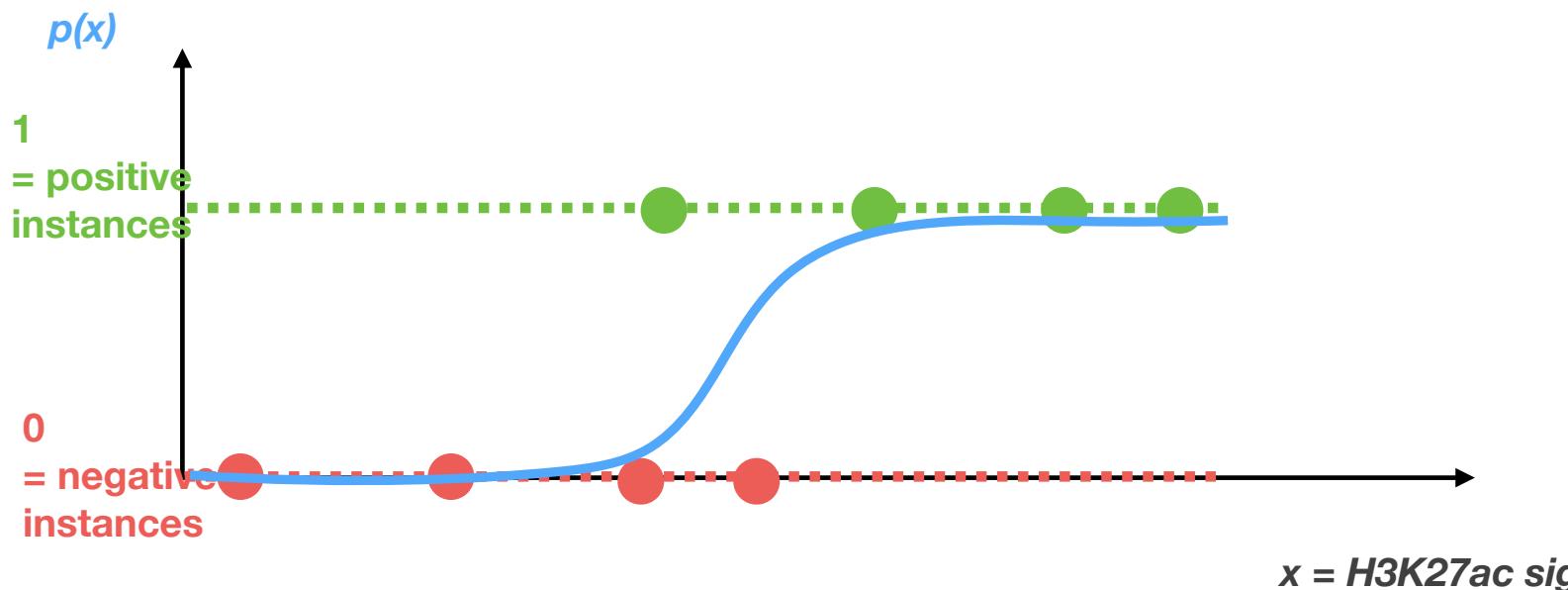
chromatin
properties

Table 1. The 29 features that were used to compute a GBP probability of transcription factor binding at specific locations

Feature no.	Feature description	Reference
1	PhastCons score for 28-way vertebrate alignment; 0 if not available	Siepel et al. 2005; Miller et al. 2007
2	PhastCons score for placental mammal subset (18 species); 0 if not available	Siepel et al. 2005; Miller et al. 2007
3	1 if PhastCons vertebrate score is available and the score is 0; 0 otherwise	Siepel et al. 2005; Miller et al. 2007
4	1 if PhastCons placental mammal score is available and the score is 0; 0 otherwise	Siepel et al. 2005; Miller et al. 2007
5	1 if PhastCons score is not available; 0 otherwise	Siepel et al. 2005; Miller et al. 2007
6	1 if part of PhastCons highly conserved vertebrate element; 0 otherwise	Siepel et al. 2005; Miller et al. 2007
7	1 if part of PhastCons highly conserved placental mammal element; 0 otherwise	Siepel et al. 2005; Miller et al. 2007
8	1 if part of a conserved indel region; 0 otherwise	Lunter et al. 2006
9	$\ln(x + 5)$ where x is distance in base pairs to nearest base of a vertebrate PhastCons element (x is 0 if base is in a highly conserved element)	Siepel et al. 2005; Miller et al. 2007
10	$\ln(x + 5)$ where x is distance in base pairs to nearest base of a placental mammal PhastCons element (x is 0 if base is in a highly conserved element)	Siepel et al. 2005; Miller et al. 2007
11	$\ln(x + 5)$ where x is distance in base pairs to nearest of a conserved indel region (x is 0 if base is in a highly conserved element)	Lunter et al. 2006
12	The estimated melting temperature at the base	Liu et al. 2007
13	Percentage of G or C base pairs of all bases within 50 bases in either direction	Karolchik et al. 2008
14	1 if base is in a UCSC Genome Browser table of CpG Islands; 0 otherwise	Karolchik et al. 2008
15	1 if base is part of a repeat element based on RepeatMasker and Tandem Repeats Finder as provided by UCSC Genome Browser	http://www.repeatmasker.org/ ; Benson 1999; Kent et al. 2002
16	1 if base is part of a transcribed region of a RefSeq gene; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
17	1 if base is between the start and end of the coding region of a RefSeq gene; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
18	1 if base is part of a RefSeq exon; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
19	1 if base is part of a RefSeq exon and within the coding region of the gene; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
20	1 if base is part of a RefSeq intron; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
21	1 if in a RefSeq 3' UTR; 0 otherwise	Karolchik et al. 2008; Pruitt et al. 2007
22	1 if in a RefSeq 5' UTR; 0 otherwise	Pruitt et al. 2007; Karolchik et al. 2008
23	$\ln(x + 5)$, where x is the absolute number of base pairs to nearest RefSeq transcription start site	Pruitt et al. 2007; Karolchik et al. 2008
24	1 if base is in a reported DNase I hypersensitive region; 0 otherwise	Boyle et al. 2008
25	$\ln(x + 1)$, where x is the number of sequence reads for the interval of the base in the summary file for CTCF	Barski et al. 2007
26	$\ln(x + 1)$, where x is the number of sequence reads for the interval of the base in the summary file for histone variant H2A.Z	Barski et al. 2007
27	$\ln(x + 1)$, where x is the sum of the number of sequence reads for the interval of the base in the summary files for the 20 histone methylation modifications.	Barski et al. 2007
28	The sum over $\ln(x_i + 1)$ for $i = 1, \dots, 20$, where the x_i 's are the number of sequence reads for the interval of the base for the 20 histone methylation modifications.	Barski et al. 2007
29	$\ln(x + 1)$, where x is the number of sequence reads in the interval of the base in the summary file for RNA polymerase II	Barski et al. 2007

Logistic regression

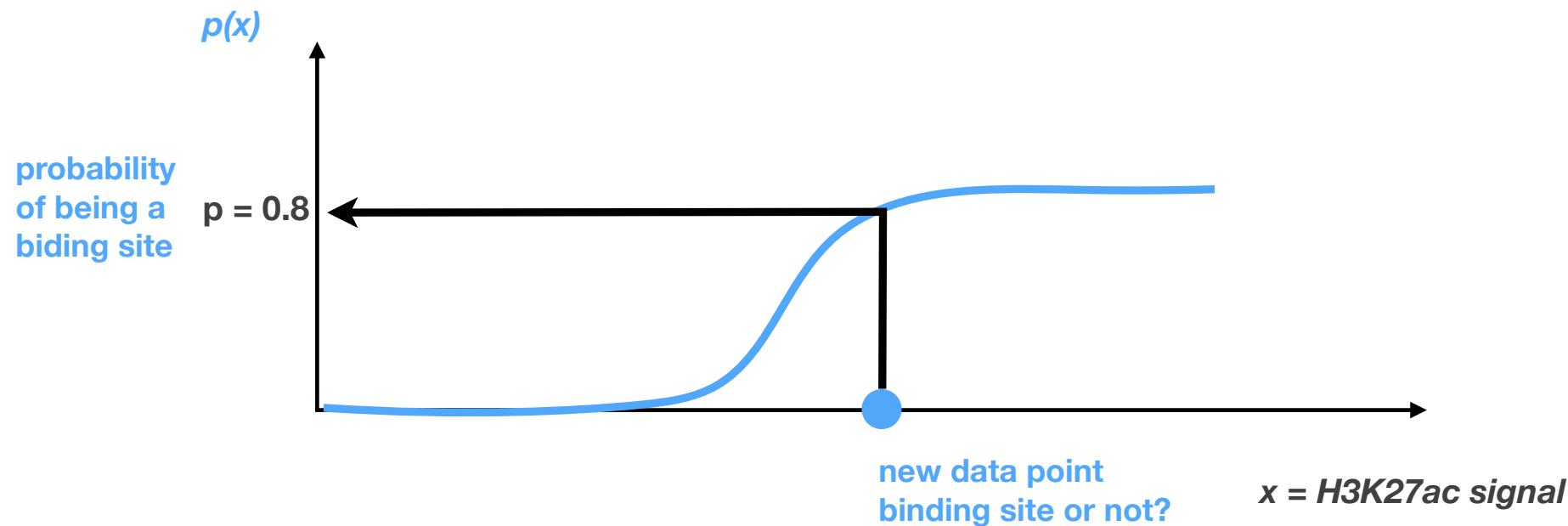
Phase 1 (training phase) : learn the model parameters from training set of positive and negative instances



$$z = \beta_0 + \beta_1 x \quad p(x) = \frac{1}{1 + e^{-z}} \quad p(x) \text{ is the } \mathbf{logit} \text{ function}$$

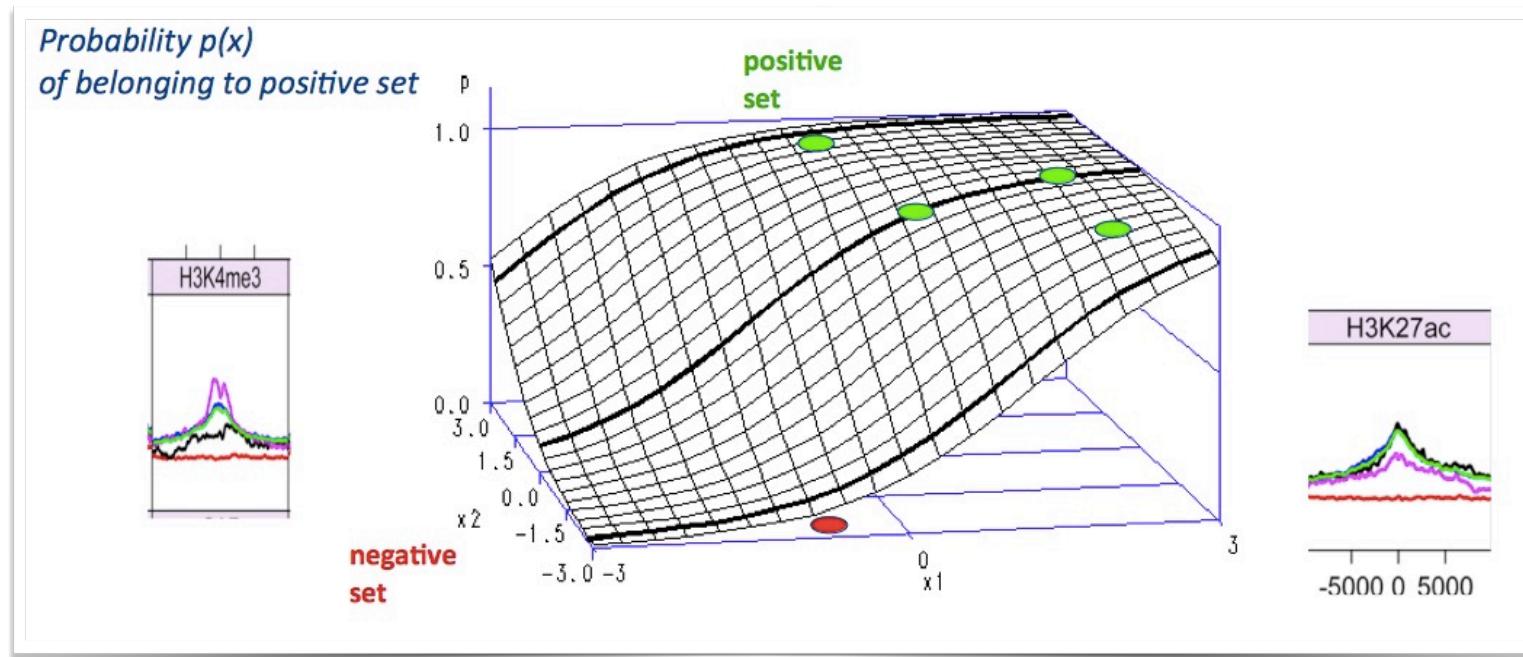
Logistic regression

Phase 2 (prediction) : predict binding probability of a new instance



Logistic regression

Using two variables (H3K4me3 and H3K27ac)



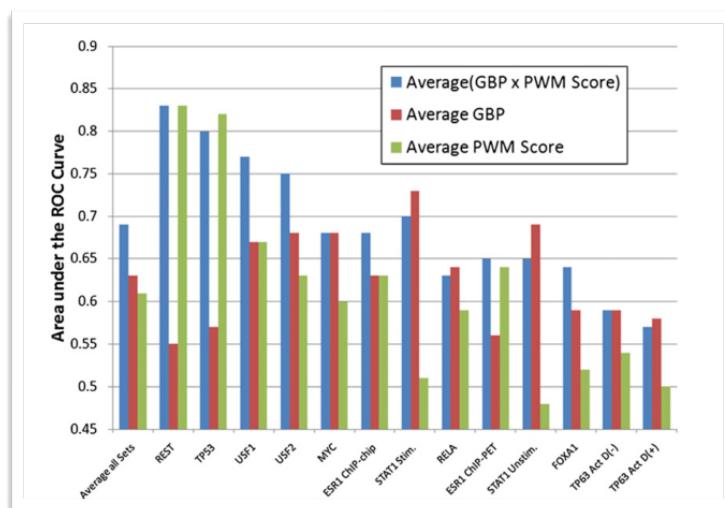
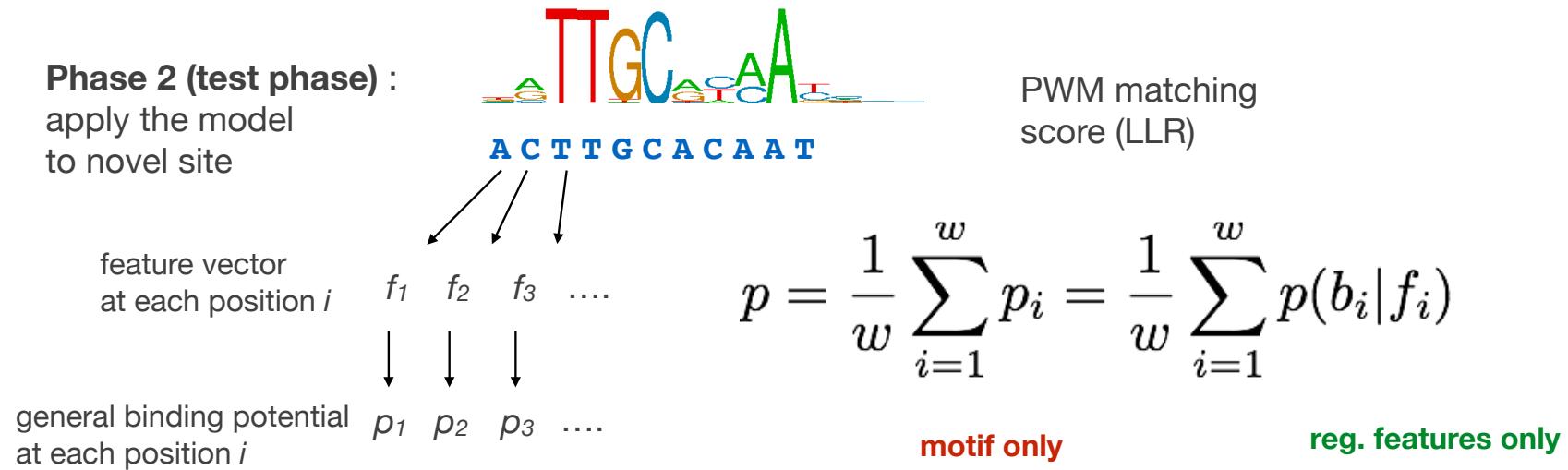
For n features :

$$z(x) = \beta_0 + \sum_{i=1}^n \beta_i x_i \quad p(x) = \frac{1}{1 + e^{-z(x)}}$$

For each nucleotide b , determine the feature vector $f = \{x_i\}$
→ Nucleotide level GBP : $p(b|f) \quad f = \{x_i\}$

General binding potential

Phase 2 (test phase) :
apply the model
to novel site



*Integrating general binding potential
to PWM score achieves better predictive
power for most motifs !*

$$S = LLR \times \frac{1}{w} \sum_{i=1}^w p(b_i | f_i)$$