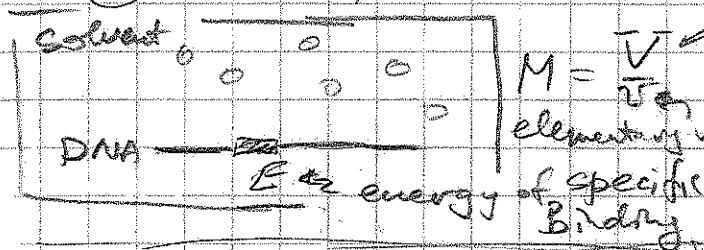


Protein-DNA interactions

① One specific site



$$Y = \frac{e^{-E} \binom{M}{P-1}}{e^{-E} \binom{M}{P-1} + \binom{M}{P}}$$

$$= \frac{e^{-E}}{e^{-E} + \binom{M}{P} / \binom{M}{P-1}} \approx \frac{e^{-E}}{e^{-E} + \frac{M}{P}}$$

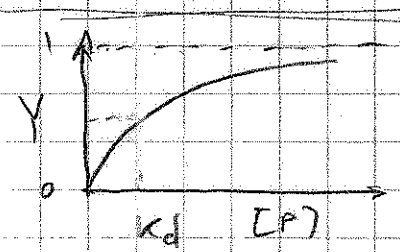
$$= \frac{P/M}{P/M + e^E} = \frac{P/V}{P/V + e^{E/V}}$$

$$= \frac{[P]}{[P] + e^{E/V}}$$

$$K_d = \frac{e^{E/V}}{V}$$

$$\frac{\binom{M}{P}}{\binom{M}{P-1}} = \frac{M! (M-P+1)! (P-1)!}{P! (M-P)! M!}$$

$$= \frac{M-P+1}{P} \approx \frac{M}{P}$$



② One specific sites, many non-specific sites

E - specific energy  
 E - non-specific energy

D - # of non-specific sites on DNA  
 = (DNA volume) / v

$$Y = \frac{e^{-E} Z(P-1)}{e^{-E} Z(P-1) + Z(P)}$$

where Z(P) - partition function for P proteins to be in the solvent or on non-specific DNA.

$$Z(P) = \sum_{k=0}^P e^{-kE} \binom{M-D}{P-k} \binom{D}{k}$$

← proteins on DNA  
 ← proteins in solvent  
 ← # of proteins on DNA

Note:

$$\binom{A}{b} = \frac{A!}{(A-b)! b!} \approx \frac{A^b}{b!}$$

b << A

then  $\binom{M-D}{P-k} \binom{D}{k} \approx \binom{M-D}{P-k} D^k \frac{1}{k! (P-k)!}$

$$= M^P \left(1 - \frac{D}{M}\right)^{P-k} \frac{D^k}{k! (P-k)!} \frac{1}{P!}$$

$$= \left(1 - \frac{D}{M}\right)^{P-k} \left(\frac{D}{M}\right)^k \binom{P}{k} \frac{M^P}{P!}$$

$$Z(P) = \left[ \sum_{k=0}^P \binom{P}{k} (1-\pi)^{P-k} \pi^k \cdot e^{-\epsilon k} \right] \frac{M^P}{P!}$$

(2)

Binomial prob of putting  $k$  particles on DNA

if  $\pi$  - prob to be on DNA;  $\pi \equiv \frac{D}{M}$

$$= (1-\pi + \pi e^{-\epsilon})^P \cdot \frac{M^P}{P!}$$

$$\frac{Z(P)}{Z(P-1)} = \left( 1 + \pi(e^{-\epsilon} - 1) \right) \frac{M}{P} \quad ; \quad \text{if } \epsilon=0 \text{ (no affinity for non-spec DNA), then } \frac{Z(P)}{Z(P-1)} = \frac{M}{P} \text{ as above.}$$

$$= \frac{M}{P} \left( 1 + \frac{D}{M} + \frac{D}{M} e^{-\epsilon} \right)$$

$$= \frac{M-D}{P} + \frac{D}{P} e^{-\epsilon}$$

$$Y = \frac{e^{-\epsilon}}{e^{-\epsilon} + \frac{Z(P)}{Z(P-1)}} = \frac{e^{-\epsilon}}{e^{-\epsilon} + \frac{M-D}{P} + \frac{D}{P} e^{-\epsilon}}$$

# of states in solvent per protein  
# of states on DNA per protein

$$= \frac{P/(M-D)}{P/(M-D) + e^{\epsilon} \left( 1 + \frac{D}{M-D} \cdot e^{-\epsilon} \right)} \cdot [P] \nu$$

recall

$$M = \frac{V}{\nu} ; D = \frac{\text{DNA vol}}{\nu}$$

$$\Rightarrow M-D = \frac{\text{solvent volume}}{\nu}$$

$$[P] \nu \cdot \frac{e^{\epsilon}}{e^{\epsilon} + \left( 1 + \frac{[DNA] \cdot \nu \cdot e^{-\epsilon}}{e^{\epsilon} / \nu} \right)}$$

where  $[P] = \frac{P}{\text{solvent volume}}$ ;  $[DNA] = \frac{D}{\text{solvent volume}}$

$$= \frac{[P]}{[P] + \frac{e^{\epsilon}}{\nu} \left( 1 + \frac{[DNA]}{e^{\epsilon}/\nu} \right)} = \frac{[P]}{[P] + K_d \left( 1 + \frac{[DNA]}{K_d^{ns}} \right)}$$

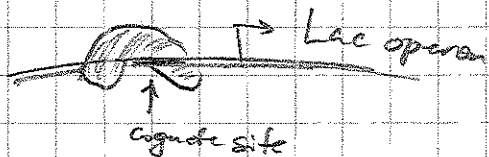
$$K_d \equiv \frac{e^{\epsilon}}{\nu} ; K_d^{ns} \equiv \frac{e^{\epsilon}}{\nu}$$

$K_d^{eff}$  - higher than  $K_d$  due to non-spec. binding

# Example: Lac repressor

3

Function:  
repress Lac operon



Lactose	Lac repressor	Lac operon	$K_d$
+	unbound	expressed	$10^{-8} M$
-	bound	repressed	$10^{-12} M$

$[P] = 10$  molecules per E. coli cell  $\approx 10 \cdot 10^{-9} M$

① Naive approach, i.e. disregarding non-specific binding

$$Y^{-lac} = \frac{[P]}{[P] + k_d} = \frac{10 \cdot 10^{-9}}{10 \cdot 10^{-9} + 10^{-12}} = \frac{1}{1 + 10^{-3}} = 0.999 \dots \text{bound!}$$

$k_d^{-lac} = 10^{-12} M$

$$Y^{+lac} = \frac{10 \cdot 10^{-9}}{10 \cdot 10^{-9} + 10^{-8}} = \frac{10}{11} = 0.9 \text{ also bound, but shouldn't}$$

② Take into account non-specific binding to DNA

$K_d^{ns} \approx 10^{-6} M$  irrespective of lactose

$[DNA] = 5 \cdot 10^6$  bp / cell /  $\sim 10$  bp footprint

$$\left(1 + \frac{[DNA]}{K_d^{ns}}\right) = 1 + \frac{5 \cdot 10^6 \cdot 10^{-9}}{10^{-6} \cdot 10} = 1 + 500 \approx 500$$

$$Y^{-lac} = \frac{10^{-8}}{10^{-8} + 10^{-12} \cdot 500} = 0.95 \text{ bound}$$

$$Y^{+lac} = \frac{10^{-8}}{10^{-8} + 10^{-9} \cdot 500} \approx 10^{-2} \text{ non-bound}$$

\* Model of binding: PWM: position weight matrix

① Assume energy is sum of contributions of bps.

$$E_{i\alpha} = \sum_{i=1}^L E_i \alpha_i \quad ; \quad \text{where } \alpha_i = \{1, 4, 1, 2, 3, \dots\}; i=1 \dots L$$

base-pair sequence of DNA  
in the site of length L

$E_{i\alpha}$  -  $L \times 4$  matrix of energy (aka PWM)

Then

$$Z = \prod_{i=1}^L \sum_{\alpha_i=1}^4 [e^{-\beta E_{i\alpha}} \cdot p_0(\alpha)]$$

background prob. of bp  $\alpha$   
for individual positions

$$P_{\alpha}(i) = \frac{e^{-\beta E_{i\alpha}} \cdot p_0(\alpha)}{\sum_{\alpha'} e^{-\beta E_{i\alpha'}} \cdot p_0(\alpha')} \leftarrow \text{freq. in SELEX experiment}$$

Const

$$\tilde{\beta E}_{i\alpha} = -\log \frac{P_i(\alpha)}{p_0(\alpha)} + \text{const};$$

Interestingly mean energy:  $\langle E \rangle = -\sum_{i,\alpha} P_i(\alpha) \log \frac{P_i(\alpha)}{p_0(\alpha)}$   
 $= -I$  information content of the motif  
 (if  $p_0(\alpha) = 1/4$ )

② Simple model of specific binding

$$E_{i\alpha} = \begin{cases} 0 & \text{for } \alpha_i^* \text{ in native site} \\ E > 0 & \text{otherwise} \end{cases}$$

then we can ask what value of  $E$  is sufficient for specific recognition of the native site of length  $L$

$$Z = N \cdot \sum_{m=0}^L \binom{L}{m} \left(\frac{1}{4}\right)^{L-m} \left(\frac{3}{4}\right)^m e^{-\beta m E}$$

↑ ~~m~~ mismatches

$$\beta F_{\text{genom}} = -\log Z$$

$$F_{\text{native}} = E - TS_{i=0} = 0$$

↑ = 0

Need  $\beta F_{\text{genomic}} \geq 0 \Rightarrow Z \leq 1$

(5)

$$Z = N \left( \frac{1}{4} + \frac{3}{4} e^{-\beta E} \right)^L \leq 1$$

$$\log_2(1 + 3e^{-\beta E}) \leq \frac{\log_2 N}{L} + 2$$

Bacteria  $N = 10^7$ ;  $L = 12$   $E \geq 1.6 kT$

mammals  $N = 10^9$   $L = 6$   $E \geq 4 kT$  !  
impossible!

⊗ Comparing to non-specific DNA

$$Z_{sp} = N \left( \frac{1}{4} + \frac{3}{4} e^{-\beta E} \right)^L \leq Z_{ns} = N \cdot e^{-\beta E_{ns}} \leq 1$$

$$\frac{E_{ns}}{L} = kT$$

⊗ Information-theoretic argument

Consider "motif" given by the frequencies  $p_i(\alpha)$   
The information content of each base-pair is

$$I_i = - \sum_{\alpha} p_i(\alpha) \log_2 \frac{p_i(\alpha)}{p_0(\alpha)} \quad (\text{in Bits})$$

Example: 1)  $p_0(\alpha) = 1/4$  for  $\alpha = A, T, G, C$ ;  $p_i(A) = 1$ ,  $p_i(G) = p_i(C) = p_i(T) = 0$   
then  $I_i = 2$  bits

2) if  $p_i(A) = p_i(G) = 1/2$ , then  $I_i = 1$  bit of information

Total information content of motif

$$I = \sum_{i=1}^L I_i = \sum_{i=1}^L \sum_{\alpha=1}^4 p_i(\alpha) \log_2 \frac{p_i(\alpha)}{p_0(\alpha)} = \langle E \rangle \quad (\text{see above})$$

computed  
 $E_{i\alpha} = - \log_2 \frac{p_i(\alpha)}{p_0(\alpha)}$

The meaning of the information content:

To make a choice out of two alternatives you need 1 bit of information;  
out of 4 alternatives: 2 Bits,  
out of N you need  $\log_2 N$  Bits.

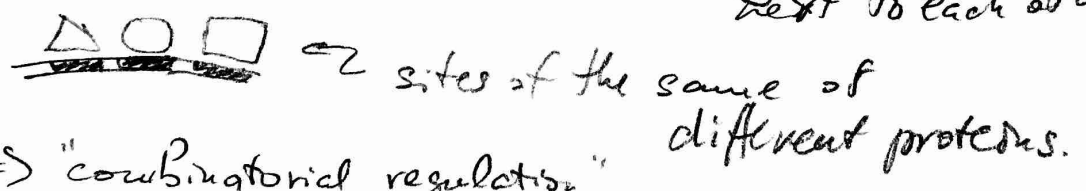
For a protein to find its site out of N alternatives in the genome of N bp, the protein needs a "antif" with  $I_{min} = \log_2 N$

Bacteria:  $I_{min} = 23$  Bits  
 $N \approx 5 \cdot 10^6 - 10^7$  Observed for bacterial proteins  $I = 23$  bits!  
(Wunderlich & Mirny, 2009)

Human  
 $N = 3 \cdot 10^9$   $I_{min} \approx 31$  bit  
 $I_{observed} \approx 12$  bits  $\ll I_{min}$ !

Even when only accessible 1% of DNA is considered  
 $N = 3 \cdot 10^7$   $I_{min} = 25$  bits  $\gg 12$  bits!

Sufficient information is provided by 2-3 sites next to each other



$\Rightarrow$  "combinatorial regulation"  
a single binding of protein to DNA is not sufficient as it happens  $\sim$  every  $2^I \approx 4000$  pb  
( $I = 12$  bits)