- Protein – DNA interactions
  - Functional: programmable : $k_d$ ; $k_d^{nS}$
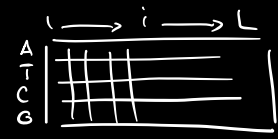
  - Specificity : $E = \sum_{i=1}^{L} \mathcal{E}(i, b_i)$

  

  $(a)$ inferred from experimental data
  $(b)$ genome-wide binding

  $i = 1, 2, \dots, L$

  - kinetics

○ Inference of $\mathcal{E}(i, b)$

Berg & von Hippel

$i \longrightarrow i \longrightarrow L$



from the data

collection of sequences bound by the protein

○ if positions were independent:
$$f(i, b) = e^{-\mathcal{E}(i,b)} / \sum_a e^{-\mathcal{E}(i,a)}$$

frequency of B at position i

$\Rightarrow \mathcal{E}(i, b) = -\log f^{obs}(i, b) + Const$
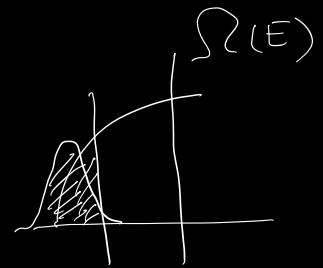
$= +\log \dfrac{f^{obs}(i, \emptyset)}{f^{obs}(i, b)}$

the most frequent at position i

○ non independence of contributions from positions
$$f(i, b) = e^{-\lambda \mathcal{E}(i,b)} / \sum_a e^{-\lambda \mathcal{E}(i,a)}$$

$E^{obs}\begin{bmatrix} selected \\ sequences \end{bmatrix} = \overline{E}$ ⟵ ?

$\Omega(E)$



$\underline{\lambda \mathcal{E}(i, b)} = \log \dfrac{f^{obs}(i, \emptyset)}{f^{obs}(i, b)}$

○ 2006 and later
Inference of $\underline{\mathcal{E}(i, b)}$

to max likelihood of observed sequences

Note: Experiments:



random

expression

Berg and von Hippel
Selection of DNA binding sites by regulatory proteins.
Statistical-mechanical theory and application to operators and promoters
https://pubmed.ncbi.nlm.nih.gov/3612791/

(b) Binding in the genome-wide context
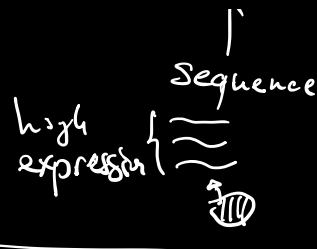
- genome is random

$$p(A) = p(T) = p(C) = p(G) = 1/4$$

- $M$ - length of genome
- Bacteria
- NO non-specific binding
- cognate site $S^*$ ← minimal energy $E(s^*) = 0$

$$E_{ib} = \begin{matrix} A \\ T \\ C \\ G \end{matrix} \begin{bmatrix} 0 & \varepsilon & \varepsilon & \varepsilon \\ \varepsilon & 0 & 0 & \varepsilon \\ \varepsilon & \varepsilon & \varepsilon & 0 \\ \varepsilon & \varepsilon & \varepsilon & 0 & 0 \end{bmatrix} \quad ; \quad \varepsilon > 0$$

→ Probability of $S^*$ occupied

$$S^* = \{ATTCGG\}$$

$$P(s^*) = \frac{Z_*}{Z_* + Z_{sp}}$$

$$Z_* = e^{-E(s^*)} \quad ; \quad \uparrow = 1$$

$$Z_{sp} = M \cdot \sum_{k=0}^{L} \left(\frac{3}{4}\right)^k \left(\frac{1}{4}\right)^{L-k} \binom{L}{k} \cdot e^{-\varepsilon \cdot k}$$

$\underset{k\text{-position mutants}}{\underbrace{\phantom{xxxxxx}}}$ ← E.g. $k = 0$

$$M \cdot \left(\frac{1}{4}\right)^L \quad \leftarrow \begin{array}{l} \text{exact} \\ \text{matches} \\ \text{of } S^* \\ \text{in the random} \\ \text{genome} \end{array}$$

$$= M \left(\frac{1}{4} + \frac{3}{4} e^{-\varepsilon}\right)^L$$

under what conditions

$$Z_{sp} \leq 1$$

$$M \left(\frac{1}{4} + \frac{3}{4} e^{-\varepsilon}\right)^L \leq 1$$

$$Z_{sp} = M \gg 1$$

E.g. $\varepsilon = 0 \implies P(s^*) = \frac{1}{M}$

$$\varepsilon \to \infty$$

$$Z_{sp} = M \frac{1}{4^L} \qquad \underline{Z_* = 1}$$
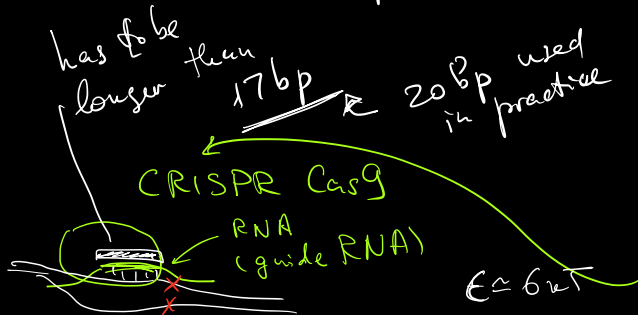
$$M / 4^L \leq 1 \qquad 4^L \geq M$$

$$\boxed{L_{crit} \geq \frac{\log_2 M}{2}}$$

For bacteria:
$M = 5 \cdot 10^6$ ; $\log_2 M = 22$ Bits $\boxed{L_{crit} = 11\ bp}$

For human:
$M = 3 \cdot 10^9$ ; $\log_2 = 32$ Bits ; $\boxed{L_{crit} = 16\ bp}$

has to be then
longer $17\ bp \approx 20\ bp$ used in practice

CRISPR Cas9

RNA (guide RNA)

$\varepsilon \approx 6\ kT$

high expression $\{$ ≈ ≈ ≈ Sequence

* Estimate $\varepsilon$ for Bacteria $\quad L = 15-20 \, bp > L_{crit}$

$$\frac{M}{4^L}\left(1 + 3e^{-\varepsilon}\right)^L \le 1$$

$$1 + 3e^{-\varepsilon} \le 4 M^{-1/4}$$

$$\varepsilon \ge \log\left[\frac{4 M^{-1/4} - 1}{3}\right]$$

$$\boxed{\varepsilon \ge 2 kT} \qquad L = 15 \, bp$$

H-Bond 1-2

or small
surface
for hydrophobic
interactions



* How about non-specific binding

$$E_{ns} = \varepsilon_{ns} \cdot L \ge 0 \qquad P(s^*) = \frac{Z_*}{Z_* + Z_{sp} + Z_{ns}}$$
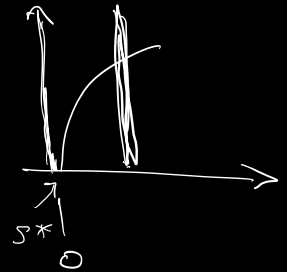
$$Z_{ns} \le 1$$

$$Z_{ns} = M \cdot e^{-E_{ns}} = M \cdot e^{-\varepsilon_{ns} \cdot L} \le 1$$

$$\varepsilon_{ns} \ge \frac{\log M}{L} \quad ; \text{ for Bacteria } L = 15 \, bp$$

$$\boxed{\varepsilon_{ns} \ge 1 kT}$$

$E_{ns}$



$s^*$ $\quad 0$

* Information theory approach

bacteria: 22 bits

$$I_{required} = \log_2 M$$

$\qquad \qquad \qquad$ human $32$ bits

$M$

$\varepsilon_{ib}$ $\qquad \varepsilon \to \infty$ $\qquad$ motif

$$\text{A T G G} \qquad I = 2 \cdot L$$

2 bits $\quad$ 2 bits $\quad$ L

$$2L \ge \log_2 M \quad ; \quad \boxed{L \ge \frac{\log_2 M}{2}}$$

$t = finite$

C    T    G

A    T    A

$I = 1 bit + 2 bit + 1 bit = 4 bit$

$I^{experimental}_{motifs} \quad \underrightarrow{} \quad I^{required}$

Different gene regulation strategies revealed by analysis of binding motifs
https://pubmed.ncbi.nlm.nih.gov/19815308/