# Protein-DNA interactions

1. For regulation ⟹ ...
2. Specific recognition : estimates
3. kinetics: protein-DNA search

---

1. Lac repressor [LacI] regulates Lac operon

E. coli    glucose - prefered

if lactose is present ⟹ Lac operon active

Lac operon

binding site of **LacI**

Binding Site / Lactose Site

Lac operon $K_d$

| | not bound | | |
|---|---|---|---|
| + | | + | $10^{-9}$ M |
| − | bound + | − | $10^{-12}$ M |

↑
occupancy of the binding site on DNA

$$Y = \frac{[P \cdot S]}{[P \cdot S] + [S]}$$

↗ occupance

$$P + S \underset{}{\overset{K_d}{\rightleftharpoons}} P \cdot S \quad ; \quad K_d = \frac{1}{v} e^{-\beta E_b} [M]$$

↳ standard volume
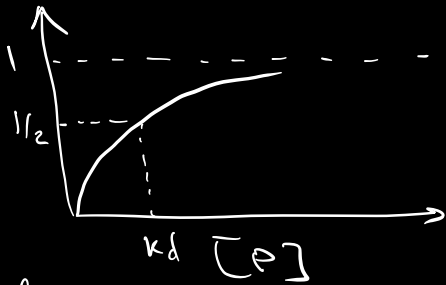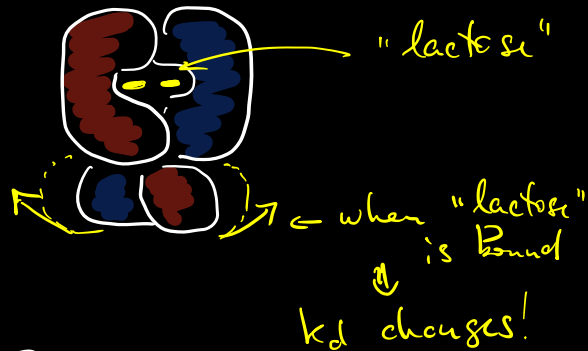
$$K_d = \frac{[P][S]}{[P \cdot S]}$$

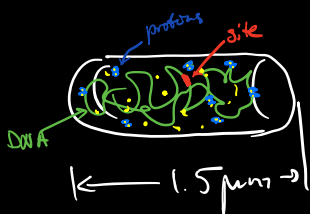$$Y = \frac{[P\cdot S]/[P][S]}{[P\cdot S]/[P][S] + 1/[P]} = \frac{1/k_d}{1/k_d + 1/[P]} = \frac{[P]}{[P] + k_d}$$



$K_d \longleftarrow$ depends on lactose

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

$K_d^{+lac} = 10^{-9} M$

$\Rightarrow$ let's compute $Y^{+lac}$

"lactose"

$\longleftarrow$ when "lactose" is bound

$\Downarrow$

$k_d$ changes!

$[P] \approx 10$ molecules per Bacteria

$\curvearrowright$ "dimers" (very stable!)

protons    site

DNA

$\updownarrow$ 1 $\mu m$    $V \simeq 1.5 \, \mu m^3$

$\longleftarrow$ 1.5 $\mu m$ $\longrightarrow$

Human cells nucleus $\sim 500 \, \mu m^3$

1 Mole = $\frac{1 \text{ molecule}}{\text{liter}}$

$\frac{1 \text{ molecule}}{1.5 \, \mu m^3} \longrightarrow \frac{1 \text{ mol}}{\text{liter}}$

10 cm

$10^{15} \, \mu m^3$

$$\boxed{\frac{1 \text{ molecule}}{1.5 \, \mu m^3} = 10^{15} / \underbrace{1.5 \cdot 6 \cdot 10^{23}}_{10} = 10^{-9} M}$$
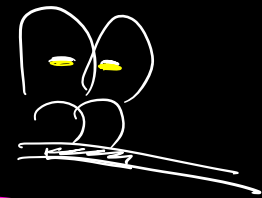
$[P] = 10 \frac{\text{molecules}}{\text{Bacteria}} = 10^{-8} M$

$$Y^{-lac} = \frac{10^{-8}}{10^{-8} + 10^{-12}} \simeq 1 \longleftarrow$$

Bound    off

$$Y^{+lac} = \frac{10^{-8}}{10^{-8} + 10^{-9}} \simeq \frac{10}{11} \simeq 0.9 \longleftarrow$$

$\longrightarrow$ should be unbound    Bound

— protein... dimer?

— kinetics... ← fast...

— DNA (non-specific DNA)

$$P + DNA \underset{ns}{\rightleftarrows} P \cdot DNA$$

$$K_d^{ns} = \frac{1}{v} e^{-\beta E_{ns}} \quad ; \quad K_d^{ns} \approx 10^{-5} M$$

$$K_d^{ns} = \frac{[P][DNA]}{[P \cdot DNA]}$$

$$[P]_{TOT} = [P] + [P \cdot DNA] + [P \cdot S]$$

$$[P] - ?$$

$$K_c^{ns} [P \cdot DNA] = [P][DNA]$$
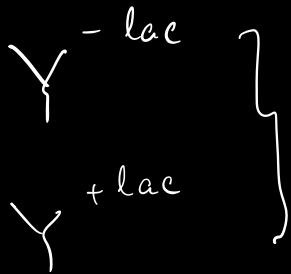
$$\frac{K_d^{ns}}{[DNA]} \left([P]_{TOT} - [P]\right) = [P]$$

$$[P] = [P]_{TOT} \Bigg/ \left(1 + \frac{[DNA]}{K_d^{ns}}\right)$$

$$Y = \frac{[P]}{[P] + K_d} = \frac{[P]_{TOT}}{[P]_{TOT} + K_d\left(1 + \frac{[DNA]}{K_d^{ns}}\right)}$$

↗ occupancy of the site

non specific
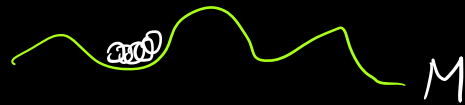
$$K_d^{eff}$$

$$Y^{-lac} \Bigg\} \quad K_d^{ns} = 10^{-6} M$$

$$Y^{+lac}$$

$$[DNA] =$$

$$= 5 \cdot 10^6 \cdot 10^{-9}$$

$$= 5 \cdot 10^{-3} M$$

$$\left(1 + \frac{[DNA]}{K_d^{ns}}\right) = 1 + \frac{5 \cdot 10^{-3}}{10^{-5}} \approx 500$$

$$M$$

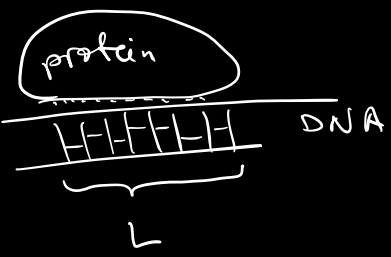$$M = 5 \cdot 10^6 \, Bp$$

↑ sites for ns binding

$$Y^{-lac} = \frac{10^{-8}}{10^{-8} + 10^{-12} \cdot 500} = \frac{10^{-8}}{10^{-8} + 10^{-10} \cdot 5} \simeq 1 \quad \boxed{Bound!}$$

$$Y^{+lac} = \frac{10^{-8}}{10^{-8} + 10^{-9} \cdot 500} = \frac{10}{10 + 500} = \underline{0.02} \quad \text{not Bound !}$$

- Non-specific DNA plays an important role! in making this system "programmable"

- Protein - DNA interactions $\Rightarrow$ logic of gene ~~regulation~~

Probing Transcription Factor Dynamics at the Single-Molecule Level in a Living Cell
Johan Elf,Gene-Wei Li,and X. Sunney Xi
Science 2007
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2853898/

② Specificity

protein

DNA

L

$$E = \sum_{i=1}^{L} \mathcal{E}(i, b_i)$$

$\uparrow$ base pair at position $i$

$\mathcal{E}(i, x)$

$L \searrow \quad \searrow x$

1. How to learn $\mathcal{E}(i,x)$?

2. How does recognition work?

| | 1 | 2 | 3 | | |
|---|---|---|---|---|---|
| A | -1 | 0 | . | | |
| T | 3 | 0 | . | | |
| | 1 | -1 | . | | |
| C | -2 | 1 | . | | |
| G | | | | | |

1. Inferring $\mathcal{E}(i,x)$ $\longrightarrow$ measure in vitro

$\Rightarrow$ difficult

$\Rightarrow$ in vitro

- In vitro PBM

site

mutant

$K_d$ $K_d$ $[P]$

$L \times 4$ experiments

Protein abundence on sequences
of the array [ M. Bulyk
G. Church ]

$\Rightarrow$ list of Bound
Sequences

In vitro
• SELEX seq

Bound DNA



random DNA sequences

Sequencing $\Rightarrow$ list of Bound Sequences

• ChIP ( Cut-and-Run ) In vivo
isolate Bound



sequence $\Rightarrow$ list of Bound sequences

?
$\downarrow$

$\mathcal{E}(i, X)$

• Evolutionary (In vivo)

LacI



LacI

LacI

E. coli
V. cholera
A. pestic

, Infer $\mathcal{E}(i, X)$ from know Bound sequences

$f(i, X) \leftarrow$ frequency of $x$ at $i$

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

$\uparrow$
$f(1, A) = \frac{2}{3}$

Boltzman equilibrium in the evolution of sites
$$p_i(x) \approx e^{-\beta \mathcal{E}(i,x)} \cdot p_0(x) / Z$$
background in the genome

$$\beta \mathcal{E}(i, x) = -\log \frac{P_i(x)}{P_0(x)} + const$$

$$\boxed{\beta \mathcal{E}(i, x) = -\log \frac{f(i, x)}{P_0(x)} \Bigg] + const}$$

observed freq in the genom

1989
Berg & von Hippel

2. (2006 ...)   H. Busemaker ( Columbia Univ )
                              Justin Kinney (Princeton)

$$L(seq, \mathcal{E}(i, x)) \xrightarrow[max]{}$$

observed sequences   $\{\mathcal{E}(i, x)\}$  MCMC

$$\Longrightarrow \mathcal{E}(i, x) \text{ matrices for lots of proteins}$$

for Bacteria /Yeast

in vitro ≃ in vivo

→ to find the site
out of $10^6 - 10^9$ alternatives
→ kinetics

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

Inference methods:
1. Statistical mechanical modeling of genome-wide transcription factor occupancy data by MatrixREDUCE
https://pubmed.ncbi.nlm.nih.gov/16873464/
2. Precise physical models of protein-DNA interaction from high-throughput data
https://pubmed.ncbi.nlm.nih.gov/17197415/