

## *Berg and Purcell revisited.*<sup>1</sup>

Pankaj Mehta

April 10, 2019

In this lecture we discuss some simple aspects of gene expression. We discuss how mRNA production is quite slow and that proteins are produced in bursts. We show how these two elements give a simple prediction for steady-state gene expression profiles.

### *The basic scaling argument*

Imagine a cell of size  $a$  wants to estimate the concentration of some chemical in the environment. We know that such an estimation is limited by diffusion and thermal noise. What are the physical limits of the best anyone could do? Berg and Purcell asked and discussed this problem in their 1977 classic paper "The Physics of Chemoreception". To answer this question, they asked formulated some abstract ideas of sensing and discussed the limits that any molecular sensor had to obey. This ideas about asking the best one could do given physical constraints has a long tradition in theoretical physics (think about the Carnot engine) and this paper brought these kind of ideas to the cellular realm.

The basic argument is quite simple. Consider a cell that occupies a volume  $V = a^3$ . If the concentration of the chemical is  $c_0$ , then, we know that on average in this volume there will be about  $\bar{N}_v = c_0 V \sim c_0 a^3$  molecules that the cell must estimate the concentration from. However, due to thermal noise, we know that there will be fluctuations in this number. We can ask about what the typical fluctuations  $\delta c$  will look like. From the central limit theorem, we know that for an instantaneous measurement

$$\frac{\delta c}{\bar{c}} \sim \frac{1}{\sqrt{N}} = \frac{1}{\sqrt{cV}} \quad (1)$$

Recalling that an *E. coli* is about a  $1(\mu m)^3$  and that a concentration of  $1nM$  corresponds to 1 molecule per bacterium. We see that using an instantaneous measurement the typical error will scale like 1 or 100% for a  $1nM$  concentration and  $1/\sqrt{1000}$  or approximately 3% for  $100\mu M$  concentration. We know that cells regularly respond to 1 –  $10nM$  concentrations indicating that cell must be doing something more.

One simple thing that cells can do is take more than one measurements. We know that if we take  $M$  independent measurements then we can reduce the error by a factor of  $\sqrt{M}$ . Notice, that the key point is of course the measurements should be uncorrelated. This is

<sup>1</sup> These notes are based on Berg Purcell "Physics of Chemoreception" 1977, Endres and Wingreen "Maximum Likelihood and the Single Receptor" PRL 2009, and Kaizu et al, "The Berg-Purcell Limit Revisited", Biophysical Journal 2014 976.

necessary because there is obviously no new information gained by measuring the same thing over and over again. In the cell, we need a completely new set of molecules to diffuse into the cell volume  $V$ . We can ask what is the characteristic time that it will take for diffusion to refresh the molecules in some volume  $V \sim a^3$ . We can do this by simple dimensional analysis. Denote the diffusion constant for the molecules that make up  $c$  by  $D$ . A biologically realistic number for small molecules is  $D = 10^{-5} \text{cm}^2/\text{s}$ . Notice that this allows to make a characteristic time for turnover by diffusion

$$\tau = D^{-1}V^{2/3} = D^{-1}a^2. \quad (2)$$

Using  $a = 10^{-4} \text{cm}$  and plugging in we get that  $\tau \approx 10^{-3} \text{s}$ . Thus, for a volume the size of *E. coli* that typical time for a concentration to renew is 1ms.

If the cell measures that concentration for a time  $T$ , the number of independent measurements will be just

$$M = T/\tau = TDV^{-2/3}, \quad (3)$$

and the typical error will be

$$\frac{\delta c}{\bar{c}} \sim \frac{1}{M\sqrt{N}} = \frac{1}{\sqrt{cV}} \times \frac{1}{\sqrt{TDV^{-2/3}}} = \frac{1}{\sqrt{cTDV^{1/3}}} = \frac{1}{\sqrt{cTDa}} \quad (4)$$

Notice that in 0.01s we can reduce the error by a factor of 10 so that for a  $100\mu\text{M}$  the error become of order 0.3%. This is an incredible accuracy.

### *Thinking about receptors*

Thus far we have just been concerned with an entire sphere. However, we know that in real cells that sensing is actually done by receptors. The first fact that Berg-Purcell show is that even a moderate number of receptors scattered over the surface of a cell (occupying less than 1/1000 of the surface) are enough to really treat the cell as a completely absorbing sphere <sup>2</sup>. To show this BP and note that at steady-state the diffusion equation

$$0 = D\nabla^2 c \quad (5)$$

is the same as the electrostatic equation with  $c$  identified with the electric potential  $\phi$ . Receptors are equipotential surface. We know that we can build a Faraday cage by having sparse metal cage, the receptors essentially can do the same thing for the chemical concentration.

Inspired by this, BP go on to calculate the uncertainty in the concentration that can be estimated from a time-series of binding and

<sup>2</sup> It is worth reading this full derivation and the beautiful analogies with electorstatistics

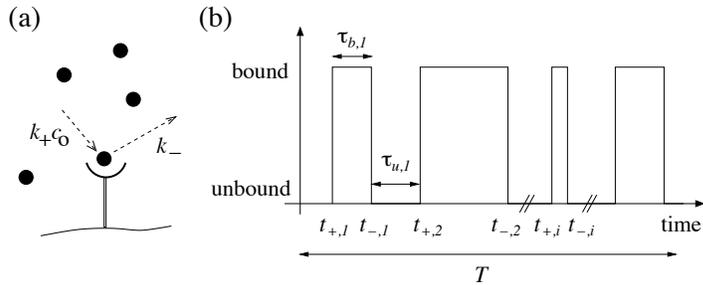


Figure 1: Figure from Endres Wingreen PRL 2009.

FIG. 1: Schematic of particle-receptor binding. (a) An unoccupied receptor can bind a particle with rate  $k_+c_0$ , and an occupied receptor can unbind a bound particle with rate  $k_-$ . (b) Binary time series of receptor occupancy.

unbinding events for a single receptor (see Figure 1). We will derive these results in a slightly different way that has much more intuition and generalizes better. Notice that to estimate the concentration we must count the number of binding events  $n$  for a duration of time  $T$ . In particular, we know that this will be a function of the concentration  $c$ :  $\bar{n}(c)$ . The cell is essentially performing an inference problem where it has to estimate  $c$  using the summary statistic  $n$ . Our error in the estimate of  $c$  comes from the fact that  $n$  will fluctuate due to thermal fluctuations. Our strategy will be to figure out the typically fluctuations in  $n$  (which we denote  $\delta n$ ) and relate these to uncertainty in our parameter estimation  $\delta c_{\text{est}}$ .

To do this, we make use of basic things we learn in all lab classes about propagating errors. We know that in the usual way that we propagate errors that

$$\delta c = d\bar{n}dc\delta c_{\text{est}} \quad (6)$$

which we can invert to get that

$$\frac{\langle \delta c_{\text{est}}^2 \rangle}{c_{\text{est}}^2} = 1/c_{\text{est}}^2 \left( \frac{dn}{dc} \right)^{-2} \langle (\delta n)^2 \rangle = 1/c_{\text{est}}^2 \left( \frac{dc}{d\bar{n}} \right)^2 \langle (\delta n)^2 \rangle \quad (7)$$

This is just the statement that we can locally linearize our sensor and the fluctuations are amplified (or dampened) by our static gain:  $\frac{dn}{dc}$ .

This leaves us with two tasks. First, we must find the function  $n(c)$ . Second, we must estimate the fluctuations in the number of binding events  $n$ . It will turn out that this is pretty easy. Let us first calculate  $n(c)$ . Let us define the mean time that the receptor is unbound as  $\bar{\tau}_u$ . Notice that this is just inversely proportional to the mean binding rate

$$\bar{\tau}_u = (k_{\text{on}}c)^{-1}. \quad (8)$$

We can also calculate the mean time that a molecule is bound before unbinding. This is just give by the inverse of the disassociation rate  $k_d$

$$\bar{\tau}_b = k_d^{-1}. \quad (9)$$

Notice that we have that the number of binding events is just

$$\bar{n} = \frac{T}{\bar{\tau}_u + \bar{\tau}_b} = \frac{T}{(k_{\text{on}}c)^{-1} + k_d^{-1}}. \quad (10)$$

We can rearrange this to get

$$c_{\text{est}} = \frac{1}{k_{\text{on}}} \cdot \frac{1}{\frac{T}{\bar{n}} - k_d^{-1}}. \quad (11)$$

We can differentiate this to get

$$\frac{dc_{\text{est}}}{d\bar{n}} = \frac{k_{\text{on}}T}{\bar{n}^2} c_{\text{est}}^2 \approx \frac{k_{\text{on}}T}{\bar{n}^2} c_0^2 = \frac{T}{k_{\text{on}}\bar{\tau}_u^2\bar{n}^2} \quad (12)$$

The second thing we need to do is estimate  $\langle(\delta n)^2\rangle$ . We will now make use of very similar trick. To quote Endres and Wingreen, “ to obtain  $\delta n$  for a fixed duration  $T$ , we note that this is proportional to the standard deviation  $\delta T$  for a fixed  $\bar{n}$  via

$$\delta n = \frac{d\bar{n}}{dT} \delta T. \quad (13)$$

Thus, we have

$$\langle(\delta n)^2\rangle = \left(\frac{d\bar{n}}{dT}\right)^2 \langle(\delta T)^2\rangle \quad (14)$$

However, we know that

$$\bar{T} = \bar{n}(\bar{\tau}_u + \bar{\tau}_b) \quad (15)$$

so that

$$\frac{d\bar{n}}{dT} = \frac{1}{\bar{\tau}_u + \bar{\tau}_b} \quad (16)$$

We also have (since we have  $n$  independent binding events hence variance is  $n$  times variance of single event)

$$\langle(\delta T)^2\rangle = \bar{n}(\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle) \quad (17)$$

This given

$$\langle(\delta n)^2\rangle = \left(\frac{\bar{n}}{T}\right)^2 \bar{n}(\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle) = \frac{\bar{n}^3}{T^2} (\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle) \quad (18)$$

We then put this together and get

$$\begin{aligned} \frac{\langle\delta c_{\text{est}}^2\rangle}{c_{\text{est}}^2} &= 1/c_{\text{est}}^2 \left(\frac{k_{\text{on}}T}{\bar{n}^2} c_0^2\right)^2 \frac{\bar{n}^3}{T^2} (\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle) \\ &= \frac{(k_{\text{on}}c_0)^2 (\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle)}{\bar{n}} \\ &= \frac{\langle(\delta\tau_u)^2\rangle + \langle(\delta\tau_b)^2\rangle}{\bar{\tau}_u^2} \cdot \frac{1}{\bar{n}} \end{aligned} \quad (19)$$

This calculation is actually quite general. You can actually convince yourself that this calculation does not depend on functional form of  $\tau_b$  and  $\tau_c$  if we just define  $k_{\text{on}} = \frac{d\tau_b}{dc}$ .

For the case of BP, we know that binding and unbinding are memoryless and described by an exponential process so that

$$\langle \tau_b^2 \rangle = \frac{1}{\bar{\tau}_b} \int_0^\infty dt t^2 e^{-t/\bar{\tau}_b} = 2\bar{\tau}_b^2 \quad (20)$$

and similarly we get

$$\langle \tau_u^2 \rangle = \frac{1}{\bar{\tau}_u} \int_0^\infty dt t^2 e^{-t/\bar{\tau}_u} = 2\bar{\tau}_u^2. \quad (21)$$

This gives  $\langle (\delta\tau_b)^2 \rangle = \bar{\tau}_b^2$  and  $\langle (\delta\tau_u)^2 \rangle = \bar{\tau}_u^2$ . Plugging this in gives

$$\frac{\langle \delta c_{\text{est}}^2 \rangle}{c_{\text{est}}^2} = \left(1 + \left(\frac{\bar{\tau}_b}{\bar{\tau}_u}\right)^2\right) \cdot \frac{1}{\bar{n}}. \quad (22)$$

Notice that from detailed balance that we have

$$\bar{p}/\bar{\tau}_b = (1 - \bar{p})/\bar{\tau}_u. \quad (23)$$

Using this we have

$$\frac{1}{\bar{n}} = \frac{\bar{\tau}_u + \bar{\tau}_b}{T} = \frac{\bar{\tau}_u}{T(1 - \bar{p})} = \frac{1}{k_{\text{on}}T(1 - \bar{p})}. \quad (24)$$

For diffusion limitation,  $k_{\text{on}} = 4Dac_0$  which is consistent with our earlier scaling. In Berg and Purcell, it was assumed that the mean bound time was equal to the mean unbound time giving

$$\frac{\langle \delta c_{\text{BP}}^2 \rangle}{c_{\text{BP}}^2} = \frac{2}{\bar{n}} = \frac{1}{2TDac_0(1 - \bar{p})}. \quad (25)$$

However, notice that this is not necessarily the best choice

### *Beyond BP: The Maximum Likelihood Solution*

If one thinks hard about it, it become clear the BP is not the best one could do. The reason is clear – notice that the fluctuations in both the bound and unbound intervals contribute to the noise. However, the bound intervals have no information about the concentration. So we should not focus on the bound intervals at all and only concentrate on the unbound intervals. This basic intuition can be formalized using Maximum Likelihood Estimation. In fact, one can show

$$\frac{1}{k_{\text{on}}c_{\text{ML}}} = T_u/n = \bar{\tau}_u \quad (26)$$

or

$$c_{\text{ML}} = \frac{1}{k_{\text{on}}\bar{\tau}_u} \quad (27)$$

To derive this, we can follow the paper of Wingreen and Endres.

### The MLE Estimator

Let us now turn this into a proper inference problem. Let us denote the observed time-series of on and off times by the set  $\{t_+, t_-\}$ . We would like to calculate the probability of observing this time series given that the molecule that binds the receptor occurs at a concentration  $c$ .

$$P(\{t_+, t_-\}; c) = \prod_i p_b(t_{+,i}, t_{-,i}) p_-(t_{-,i}) p_u(t_{-,i}, t_{+,i+1}) p_+(t_{+,i+1}), \quad (28)$$

where the probability for a particle to remain bound from  $t_{+,i}$  to  $t_{-,i}$  is

$$p_b(t_{+,i}, t_{-,i}) = p_b(t_{-,i} - t_{+,i}) = e^{-k_d(t_{-,i} - t_{+,i})} \quad (29)$$

and the probability for a receptor to remain unbound from  $t_{-,i}$  to  $t_{+,i+1}$  is

$$p_u(t_{-,i}, t_{+,i+1}) = p_u(t_{+,i+1} - t_{-,i}) = e^{-k_{\text{on}}c(t_{+,i+1} - t_{-,i})}. \quad (30)$$

Furthermore, notice that the probability to bind the molecule at time  $t_{+,i}$  is proportional to the binding rate and given by

$$p_+(t_{+,i}) \propto k_{\text{on}}c. \quad (31)$$

Similarly, the probability to unbind the molecule at time  $t_{-,i}$  is proportional to the unbinding rate and given by

$$p_-(t_{-,i}) \propto k_d. \quad (32)$$

Putting together all these expressions one has

$$P(\{t_{+,i}, t_{-,i}\}; c) \propto e^{-k_d T_b} \cdot e^{-k_{\text{on}}c T_u} \cdot k_d^n \cdot (k_{\text{on}}c)^n, \quad (33)$$

where  $n$  is the number of binding event in the time series (which we assume to be equal to the number of unbinding events) and  $T_{b(u)}$  is the total time the receptor is bound (unbound):

$$T_{b(u)} = \sum_i^n \tau_{b(u),i}. \quad (34)$$

We can also calculate the log-likelihood

$$l(\{t_{+,i}, t_{-,i}\}; c) = \log P(\{t_{+,i}, t_{-,i}\}; c) = -k_d T_b - k_{\text{on}}c T_u + n(\log k_d + \log(k_{\text{on}}c)) \quad (35)$$

The Maximum Likelihood solution  $\hat{c}_{ML}$  is found by finding the value of  $c$  that maximizes the probability of the observed time-series. This can be found by just differentiating the log-likelihood and setting it to zero:

$$-k_{\text{on}}T_u + n/\hat{c}_{ML} = 0 \quad (36)$$

Rearranging this gives that

$$\hat{c}_{ML} = \frac{n}{k_{\text{on}}T} = k_{\text{on}}^{-1} \frac{1}{\bar{\tau}_u}, \quad (37)$$

where  $\bar{\tau}_u = T/n$  is the average length of the unbound intervals. In other words, we know that we choose the parameter to match the average time of the unbound intervals

$$k_{\text{on}}\hat{c}_{ML} = \bar{\tau}_u. \quad (38)$$

Why does this result make sense? Notice that the duration of the bound intervals contain no information about the concentration since they do not depend on  $c$ . Thus, focusing on the bound intervals only adds noise to our estimate without giving us any more information. Thus, it is not surprising that the optimal MLE estimator only uses the unbound intervals!

### *Calculating Uncertainty using the Cramer-Rao Bound*

We would also like to calculate the uncertainty in our estimate. This seems like a very hard problem. However, we can invoke a classic result for statistics called the Cramer-Rao bound which state the uncertainty in the ML estimator is bounded by the second derivative of the log-likelihood function

$$\frac{(\delta c_{ML})^2}{c_{ML}^2} = - \frac{1}{c_0^2 \langle \frac{\partial^2 l(\{t_{+,i}, t_{-,i}\}; c)}{\partial c^2} \rangle} \quad (39)$$

This has an intuitive interpretation. The variance in our estimator is basically given by the curvature around the maximum. We can now just calculate this directly to get

$$- \frac{\partial^2 l(\{t_{+,i}, t_{-,i}\}; c)}{\partial c^2} = \frac{n}{c_0}, \quad (40)$$

or as expected

$$\frac{(\delta c_{ML})^2}{c_{ML}^2} = \frac{1}{n}. \quad (41)$$

Notice this corresponds exactly to the case where  $\langle (\delta \tau_b)^2 \rangle = 0$  and  $\langle (\delta \tau_u)^2 \rangle = \bar{\tau}_u$  as one would expect if the bound intervals were deterministic but the unbound intervals were set by diffusion!

### *"Renormalizing" rates to account for receptor re-binding*

Thus, far we ignored the fact that ligands are more likely to rebind the receptor after unbinding. In fact a real time series looks much more like Figure 2.

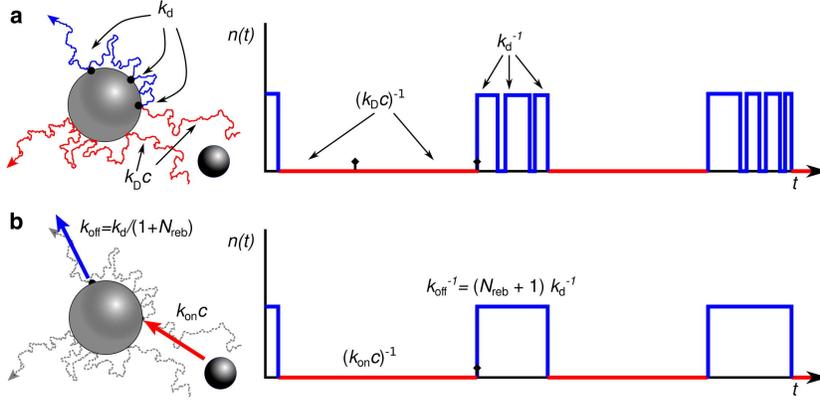


Figure 2: Actual time series has many rebinding events. However, this can be modeled by effective "renormalized" on and off rates. Figure from Kaizu et al Biophysical Journal 2014

How can we account for this. We will make use of some simple time-scale separation arguments. Before doing that, let us define some notation. Define the diffusive rate constant by  $k_d = 4\pi aD$ , where  $a$  is the cross-section of the receptor and  $D$  is the diffusion constant. Notice  $k_D$  is related to the flux of molecules incident upon the receptor with the diffusive flux given by  $k_D c$ . Let us denote by  $k_a$  the rate at which a ligand that encounters the receptor binds the receptor (Note quotes are from Kaizu et al).

- We note that the time a ligand molecule spends near the receptor is much smaller than the timescale on which ligand molecules arrive at the receptor from the bulk given by  $(k_d c)^{-1}$ .
- When a molecule unbinds from the receptor, it either rebinds with a probability  $p_{reb}$  or escapes to the bulk with probability  $1 - p_{reb}$ .
- It is clear that there is a competition between rebinding and a new molecule arriving so that  $p_{reb} = \frac{k_a}{k_D + k_a}$ . This also implies that the escape probability is just  $p_{esc} = 1 - p_{reb} = \frac{K_D}{k_a + K_D}$ .
- "The mean number of rounds of rebinding and dissociation before the molecule escapes into the bulk is then":

$$N_{reb} = (1 - p_{reb}) \sum_j j p_{reb}^j = \frac{p_{reb}}{1 - p_{reb}} = k_a / K_D \quad (42)$$

- "The average number of times a molecule from the bulk encounters the receptor before it actually binds is":

$$N_{esc} + 1 = p_{reb} \sum_j j p_{esc}^j + 1 = k_D / k_a + 1 = 1 / N_{reb} + 1 \quad (43)$$

- The effective renormalized rate at which the molecule binds is then

$$k_{on} = k_D / (1 + N_{esc}) = k_a p_{esc} = \frac{k_a K_D}{k_a + K_D} \quad (44)$$

- The effective renormalized rate at which the molecule disassociates is then

$$k_d^{eff} = k_d p_{esc} = k_d / (1 + N_{reb}) = \frac{k_d K_D}{k_a + K_D} \quad (45)$$

These heuristic arguments can be made rigorous with a real calculation<sup>3</sup>.

*Molecularly plausible realizations of these inference algorithms*

<sup>3</sup> see Kaizu et al, "The Berg-Purcell Limit Revisited", Biophysical Journal 2014 976 for the beautiful but technical calculation.