

Dynamics of double-strand break and repair in growing populations

Yeast team¹

¹*Harvard, Pasteur, Peking*

Contents

1	Introduction	1
2	Analysis of pooled data reveals inconsistencies with simple model	1
3	Fluctuation analysis	2
A	Analysis of ODE model	3
B	Analysis of asynchronous model	3
C	Analysis of synchronous model	4
C.1	Two state model	4
C.2	Broken state model	5

1 Introduction

2 Analysis of pooled data reveals inconsistencies with simple model

We begin by ignoring the single-cell resolution afforded to us by the individual wells, and treat each condition as a bulk experiment by studying the number of bright field and GFP cells averaged over all the wells. The trajectories of the average number of bright field and GFP cells are shown in Figure 2 (A). We compare this data to a simple ODE mode in which an initial population of cells with modified DNA grows at a rate α and switches into a non-growing, broken state, at a rate ρ . The broken cells can then become repaired and once again begin to grow at a rate α . Letting m , b and g be the number of modified, broken and repaired (or green) cells, we have

$$\frac{d}{dt}m = (\alpha - \beta)m \tag{1}$$

$$\frac{d}{dt}b = \beta m - \rho b \tag{2}$$

$$\frac{d}{dt}g = \alpha g + \rho b. \tag{3}$$

It is easy to see that the total population size $n = m + b + g$ will eventually grow exponentially at a rate α . We can therefore obtain α by fitting the bright field data to an exponential. A prediction of this model is that the fraction of non-GFP cells, $\phi_{m+b} = (m + b)/(m + b + g)$ decays exponentially at a rate that is independent of ρ :

$$\phi_{m+b} = 1 - e^{\beta t}. \tag{4}$$

This exponentially decay is seen in some, but not all experiments; see Figure 2 (B)

The trajectories of $r(t)$ shown in Figure 2 (C) appear inconsistent with our assumptions thus far. Therefore, we consider

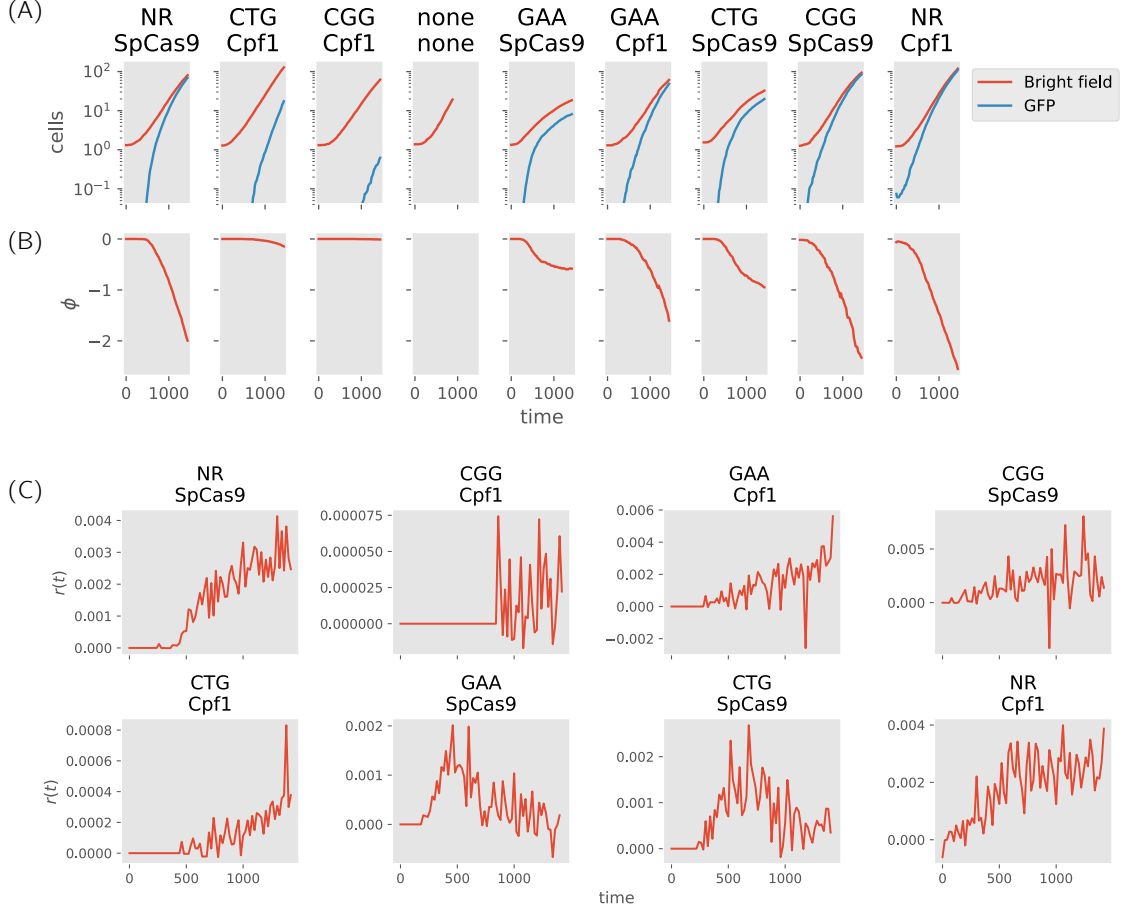


Figure 1: **Analysis of pooled data** (A) The average number of bright field and GFP cells as a function of time in each condition. (B) The fraction of non-GFP cells as function of time for each condition. (C) The rate to turn obtained from Equation ?? compared to simulations with the values of inferred from the trajectories in (B).

3 Fluctuation analysis

We now consider the variability between wells within a condition. In order to remove the effects of lag time and asynchronous growth, we consider the variance of G conditioned on the total number of cells, $\text{var}(G|N)$. In a simple model where we neglect the broken, non-green state, $\text{var}(G|N)$ undergoes phase transition at $p_c = 1/2$ (see Section C for derivation). When $p < 1/2$, $\text{var}(G|N)$ grows exponentially, while for $p > 1/2$ it decays. In all the experiments, we find that $\text{var}(G|N)$ grows exponentially.



Figure 2: (A) The average number of bright field and GFP cells as a function of time in each condition. (B) The fraction of non-GFP cells as function of time for each condition. (C) The rate to turn obtained from Equation ?? compared to simulations with the values of inferred from the trajectories in (B).

A Analysis of ODE model

B Analysis of asynchronous model

Here consider a stochastic model where divisions are asynchronous; that is, cells have different generation times. For simplicity, we will consider the case where generation times are independent and exponentially distributed. Although not biologically realistic, it captures the essential qualitative features of the dynamics, namely the growth of variance with time. We consider the probability $P(n, g, t)$ of observing n cells, g of which are green at time t and suppose that cells divide at a constant rate α , while non-green cells turn green at a rate β . Under these assumptions $P(n, g, t)$ obeys the master equation

$$\begin{aligned} \frac{d}{dt} P(n, g, t) = & \alpha(n - g - 1)P(n - 1, g, t) + \alpha(n - 1)P(n - 1, g - 1, t) \\ & + \beta(n - g + 1)P(n, g - 1, t) - [\alpha n + \beta(n - g)] P(n, g, t) \end{aligned} \quad (5)$$

For the averages, we have

$$\frac{d}{dt} \langle N \rangle = \alpha \langle N \rangle \quad (6)$$

$$\frac{d}{dt} \langle G \rangle = \alpha \langle G \rangle + \beta(\langle N \rangle - \langle G \rangle) \quad (7)$$

which implies $\langle N \rangle \sim e^{\alpha t}$ and

$$\langle G \rangle = e^{(\alpha - \beta)t} [e^{\beta t} - 1]. \quad (8)$$

We now consider the second moment:

$$\frac{d}{dt}\langle G^2 \rangle = \sum_{g,n} g^2 \frac{d}{dt} P(n, g, t) \quad (9)$$

$$= \sum_n \sum_g \alpha(n-g) g^2 P(n, g, t) + \sum_n \sum_g \alpha g(g+1)^2 P(n, g, t) \quad (10)$$

$$+ \sum_g \sum_g \beta(n-g)(g+1)^2 P(n, g, t) - \sum_n \sum_g [\alpha n + \beta(n-g)] g^2 P(n, g, t) \quad (11)$$

$$= \alpha \langle NG^2 \rangle - \alpha \langle G^3 \rangle + \alpha \langle G^3 \rangle + 2\alpha \langle G^2 \rangle + \alpha \langle G \rangle + \beta \langle NG^2 \rangle - \beta \langle G^3 \rangle \quad (12)$$

$$+ 2\beta \langle NG \rangle - 2\beta \langle G^2 \rangle + \beta \langle N \rangle - \beta \langle G \rangle - \langle NG^2 \rangle - \beta \langle NG^2 \rangle + \beta \langle G^3 \rangle \quad (13)$$

After cancelling terms, we find

$$\frac{d}{dt}\langle G^2 \rangle = 2\alpha \langle G^2 \rangle + \alpha \langle G \rangle + 2\beta \langle NG \rangle \quad (14)$$

$$- 2\beta \langle G^2 \rangle + \beta \langle N \rangle - \beta \langle G \rangle. \quad (15)$$

For the variance

$$\frac{d}{dt} \text{var}(G) = \frac{d}{dt} [\langle G^2 \rangle - \langle G \rangle^2] \quad (16)$$

$$= \frac{d}{dt} \langle G^2 \rangle - 2\langle G \rangle \frac{d}{dt} \langle G \rangle \quad (17)$$

$$= 2(\alpha - \beta) \text{var}(G) + 2\beta \text{cov}(G, N) + (\alpha - \beta) \langle G \rangle - \langle N \rangle. \quad (18)$$

We need to calculate the dynamics of the covariance to obtain a closed set of equations. Proceed in a similar manner as we did for the variance, we start with the product:

$$\frac{d}{dt} \langle GN \rangle = \alpha \langle G \rangle + 2\alpha \langle GN \rangle - \beta \langle GN \rangle + \beta \langle N^2 \rangle. \quad (19)$$

Subtracting the equations for the means gives

$$\frac{d}{dt} \text{cov}(G, N) = \frac{d}{dt} \langle GN \rangle - \langle G \rangle \frac{d}{dt} \langle N \rangle - \langle N \rangle \frac{d}{dt} \langle G \rangle \quad (20)$$

$$= \beta \text{var}(N) + (2\alpha - \beta) \text{cov}(G, N) + \alpha \langle G \rangle. \quad (21)$$

Since $\text{var}(N) \sim e^{2\alpha t}$, we conclude that $\text{cov}(N, G) \sim e^{2\alpha t}$ and therefore,

$$\text{var}(G) \sim e^{2\alpha t} \quad (22)$$

for any value of β .

C Analysis of synchronous model

Here we consider a class of synchronous model in which all growing cells have a fixed generation time.

C.1 Two state model

We start by ignoring any intermediate state the breaking of DNA and the appearance of growing GFP cells. In this case, the total number of cells at the t th generation is simply 2^t . Letting G_t be the number of green cells at time t , we have

$$G_t = 2G_{t-1} + 2 \sum_{j=1}^{2^{t-1}-G_{t-1}} \xi_{t,j} \quad (23)$$

where $\xi_{t,j}$ are independent Bernoulli random variables with probability of success p (the probability to turn green in a given generation). The first term represents the division of existing green cells, while the second represents the fact that each of the $2^{t-1} - G_{t-1}$ non-green cells in the previous generation has a probability p to turn green before dividing. The second term has a binomial distribution with mean $p(2^t - 2G_{t-1})$ and variance $p(1-p)(2^t - 2G_{t-1})$ and can therefore be approximated by a normal distribution, allowing us to rewrite Equation 24 as

$$G_t \approx 2(1-p)G_{t-1} + p2^t + \sqrt{p(1-p)(2^t - 2G_{t-1})}\eta_t \quad (24)$$

where η_t are independent standard normal random variates. For the mean (this actually be obtained directly from Equation 24 and is exact), we have

$$\langle G_t \rangle = 2(1-p)\langle G_{t-1} \rangle + p2^t = 2^t \sum_{j \leq t} (1-p)^{t-j} p \quad (25)$$

$$= 2^t [1 - (1-p)^t] \quad (26)$$

From this equation, we can deduce that the fraction of green cells approaches one for all values of p . For the variance, we make the approximation that G_{t-1} appearing in the square root can be replaced by $\langle G_{t-1} \rangle$ mean, which gives

$$G_t \approx 2(1-p)G_{t-1} + p2^t + \sqrt{2^t p(1-p)^t} \eta_t. \quad (27)$$

Taking the variance of both sides, we have

$$\text{var}(G_t) \approx 4(1-p)^2 \text{var}(G_{t-1}) + 2^t p(1-p)^t = \sum_{j \leq 2} 2^{2t-j} (1-p)^{2(t-j)} p(1-p)^j \quad (28)$$

$$= \frac{p((2-2p)^t - 1)(2-2p)^t}{1-2p} = \frac{p}{1-2p} \left[(2^t(1-p))^2 - 2^t(1-p)^t \right]. \quad (29)$$

Considering the cases $p < 1/2$ and $p > 1/2$ separately, we obtain the simple expressions:

$$\text{var}(G_t) = \begin{cases} \frac{p}{1-2p} [2(1-p)]^{2t} & \text{if } p < 1/2 \\ \frac{p}{2p-1} [2(1-p)]^t & \text{if } p > 1/2 \end{cases} \quad (30)$$

Rewriting this expression in terms of the total number of cells $N = 2^t$, we have

$$\text{var}(G_t|N) = \begin{cases} \frac{p}{1-2p} N^{2+2 \ln(1-p)/\ln(2)} & \text{if } p < 1/2 \\ \frac{p}{2p-1} N^{1+1 \ln(1-p)/\ln(2)} & \text{if } p > 1/2 \end{cases} \quad (31)$$

C.2 Broken state model

Let the probability for a modified cell to break be p and the probability for a broken cell to repair in a given generation be q . Letting B_t and G_t be the number of broken and green cells, we have

$$M_t = 2M_{t-1} - \sum_i^{2M_{t-1}} \xi_{t,\text{break},i} \quad (32)$$

$$B_t = \sum_{i=1}^{B_{t-1}} (1 - \xi_{t,\text{repair},i}) + \sum_i^{2M_{t-1}} \xi_{t,\text{break},i} \quad (33)$$

$$G_t = 2G_{t-1} + \sum_{i=1}^{B_{t-1}} \xi_{t,\text{repair},i} \quad (34)$$

where $\xi_{\text{repair},i}$ and $\xi_{\text{break},i}$ are 1 if the i th cells in the sums are broken or repaired respectively, and zero otherwise. Making a Gaussian approximation, we have

$$\begin{aligned} M_t &= 2M_{t-1}(1-p) + \sqrt{2p(1-p)M_{t-1}}\eta_{2,t} \\ B_t &= (1-q)B_{t-1} + 2pM_{t-1} - \sqrt{q(1-q)B_{t-1}}\eta_{1,t} - \sqrt{2p(1-p)M_{t-1}}\eta_{2,t} \\ G_t &= 2G_{t-1} + qB_{t-1} + \sqrt{q(1-q)B_{t-1}}\eta_{1,t}. \end{aligned} \quad (35)$$

Here, $\eta_{i,t}$ are independent standard normal random variates, but note that we have been careful to flip the signs in-front of the noise terms above to that, e.g., a positive value of $\eta_{2,t}$ in the first equation corresponds to a negative value in the second. Since we ultimately interested in the distribution of G_t conditioned on the total number of cells, $N_t = G_t + B_t + M_t$, we consider the equation of N_t :

$$N_t = 2(N_t - B_t). \quad (36)$$

Conditioned on the previous generation $\mathbf{z}_t = (M_t, B_t, G_t)^T$ has a multivariate normal distribution:

$$f_t(\mathbf{z}|\mathbf{z}') = \frac{1}{\sqrt{2\pi}\det\Sigma_t(\mathbf{z}')^{1/2}} e^{-\frac{1}{2}(\mathbf{z}-\boldsymbol{\mu}_t(\mathbf{z}'))^T \Sigma_t(\mathbf{z}')^{-1}(\mathbf{z}-\boldsymbol{\mu}_t(\mathbf{z}'))} \quad (37)$$

The mean and covariance matrix can be determined from Equation (36). For the mean, we have

$$\boldsymbol{\mu}_t = \begin{bmatrix} \langle M_t | \mathbf{z}_{t-1} \rangle \\ \langle B_t | \mathbf{z}_{t-1} \rangle \\ \langle G_t | \mathbf{z}_{t-1} \rangle \end{bmatrix} = A \mathbf{z}_{t-1} \quad (38)$$

where

$$A = \begin{bmatrix} 2-2p & 0 & 0 \\ 2p & 1-q & 0 \\ 0 & q & 2 \end{bmatrix}. \quad (39)$$

Writing the initial $\mathbf{z}_0 = (1, 0, 0)$, we find

$$\langle \mathbf{z}_t \rangle = \begin{bmatrix} \frac{(2-2p)^t}{2p((1-q)^t - (2-2p)^t)} \\ -\frac{q(2p(2^t - (1-q)^t) + (q+1)2^t((1-p)^t - 1))}{(q+1)(-2p+q+1)} \end{bmatrix}. \quad (40)$$

We find (as we could have deduced without this calculation) that the fraction of green cells approaches one as $t \rightarrow \infty$.

Now we turn to the covariance matrix. The elements of Σ_t , are determined by variance and covariance of b_t and g_t :

$$\Sigma_{t,1,1} = \text{var}(M_t | \mathbf{z}_{t-1}) = p(1-p)2M_{t-1} \quad (41)$$

$$\Sigma_{t,2,2} = \text{var}(B_t | \mathbf{z}_{t-1}) = q(1-q)B_{t-1} + p(1-p)2M_{t-1} \quad (42)$$

$$\Sigma_{t,3,3} = \text{var}(G_t | \mathbf{z}_{t-1}) = q(1-q)B_{t-1} \quad (43)$$

and

$$\Sigma_{t,1,2} = \text{cov}(M_t, B_t | \mathbf{z}_{t-1}) = 2p(1-p)M_{t-1} \quad (44)$$

$$\Sigma_{t,2,3} = \text{cov}(B_t, G_t | \mathbf{z}_{t-1}) = q(1-q)B_{t-1} \quad (45)$$

$$\Sigma_{t,3,1} = \text{cov}(M_t, G_t | \mathbf{z}_{t-1}) = 0 \quad (46)$$

The covariance matrix is therefore,

$$\Sigma_t = \begin{bmatrix} 2p(1-p)M_{t-1} & 2p(1-p)M_{t-1} & 0 \\ 2p(1-p)M_{t-1} & q(1-q)B_{t-1} + 2p(1-p)M_{t-1} & q(1-q)B_{t-1} \\ 0 & q(1-q)B_{t-1} & q(1-q)B_{t-1} \end{bmatrix}. \quad (47)$$

We are interested in the unconditional distribution $P_t(\mathbf{z})$, so we write the recursive equation

$$P_t(\mathbf{z}) = \int f_t(\mathbf{z}|\mathbf{z}') P_{t-1}(\mathbf{z}') d\mathbf{z}'. \quad (48)$$

We make ansatz that $P_t(\mathbf{z})$ is a Gaussian

$$P_t(\mathbf{z}) = \frac{1}{\sqrt{2\pi} [\det\Omega_t]^{1/2}} e^{-\frac{1}{2}(\mathbf{z}-\boldsymbol{\nu}_t)^T \Omega_t^{-1}(\mathbf{z}-\boldsymbol{\nu}_t)}. \quad (49)$$

This is equivalent to assuming that B_{t-1} and G_{t-1} appearing in the elements of the covariance matrix can be replaced by their averages, thus making the joint distribution of \mathbf{z} and \mathbf{z}' also a multivariate Gaussian:

$$P_t(\mathbf{z}, \mathbf{z}') = f_t(\mathbf{z}|\mathbf{z}')P_{t-1}(\mathbf{z}') \propto e^{-\frac{1}{2}(\mathbf{z}-\boldsymbol{\mu}_t(\mathbf{z}'))^T \Sigma_t^{-1}(\mathbf{z}-\boldsymbol{\mu}_t(\mathbf{z}')) - \frac{1}{2}(\mathbf{z}'-\boldsymbol{\nu}_{t-1})^T \Omega_{t-1}^{-1}(\mathbf{z}'-\boldsymbol{\nu}_{t-1})} \quad (50)$$

$$\propto e^{-\frac{1}{2}(\mathbf{z}-\hat{\boldsymbol{\mu}}_t)^T \hat{\Sigma}_t^{-1}(\mathbf{z}-\hat{\boldsymbol{\mu}}_t)} \quad (51)$$

where $\hat{\boldsymbol{\mu}}_t$ and $\hat{\Sigma}_t$ are the mean and covariance matrices of the joint distribution of the current and previous generations. Since \mathbf{z} and \mathbf{z}' have marginal distributions given by $P_t(\mathbf{z})$ and $P_{t-1}(\mathbf{z})$ respectively, $\hat{\Sigma}_t$ has the block form:

$$\hat{\Sigma}_t = \begin{bmatrix} \Omega_t & \hat{\Sigma}_{1,2} \\ \hat{\Sigma}_{1,2}^T & \Omega_{t-1} \end{bmatrix}. \quad (52)$$

where

$$\hat{\Sigma}_{1,2} = \begin{bmatrix} \text{cov}(M_t, M_{t-1}) & \text{cov}(M_t, B_{t-1}) & \text{cov}(M_t, G_{t-1}) \\ \text{cov}(B_t, M_{t-1}) & \text{cov}(B_t, B_{t-1}) & \text{cov}(B_t, G_{t-1}) \\ \text{cov}(G_t, M_{t-1}) & \text{cov}(G_t, B_{t-1}) & \text{cov}(G_t, G_{t-1}) \end{bmatrix}. \quad (53)$$

After some algebra, we can now write $\hat{\Sigma}_{1,2}$ in terms of the entries of Ω_{t-1} (dropping the time subscript for clarity):

$$\hat{\Sigma}_{1,2} = \begin{bmatrix} 2(1-p)\Omega_{1,1} & 2(1-p)\Omega_{1,2} & 2(1-p)\Omega_{1,3} \\ (1-q)\Omega_{2,1} + 2p\Omega_{1,1} & (1-q)\Omega_{2,2} + 2p\Omega_{1,2} & (1-q)\Omega_{2,3} + 2p\Omega_{1,3} \\ 2\Omega_{3,1} + q\Omega_{2,1} & 2\Omega_{3,2} + q\Omega_{2,2} & 2\Omega_{3,3} + q\Omega_{2,3} \end{bmatrix} \quad (54)$$

$$= D\Omega_{t-1} \quad (55)$$

where

$$D = \begin{bmatrix} 2(1-p) & 0 & 0 \\ 2p & (1-q) & 0 \\ 0 & q & 2 \end{bmatrix}. \quad (56)$$

In order to derive an equation for Ω_t , we notice that Σ_t (which is known) can be expressed as the shur complement

$$\Sigma_t = \Omega_t - \hat{\Sigma}_{t,1,2}\Omega_{t-1}^{-1}\hat{\Sigma}_{t,1,2}^T \quad (57)$$

which gives us a recursive relation for Ω_t :

$$\Omega_t = \Sigma_t + D(D\Omega_{t-1})^T = \Sigma_t + D\Omega_{t-1}D^T \quad (58)$$

Taking $\Omega_0 = 0$,

$$\Omega_t = \sum_{j \leq t} D^{t-j} \Sigma_j (D^T)^{t-j} \quad (59)$$

We are actually interested in the variation in green cells conditioned on the total number of cells, thus, we introduce the transformation

$$L = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}. \quad (60)$$

This gives us the covariance matrix for (N_t, G_t) :

$$\bar{\Omega}_t = L\Omega_t L^T = \sum_{j \leq t} L D^{t-j} \Sigma_j (D^T)^{t-j} L^T \quad (61)$$

In terms of the entries of $\bar{\Omega}_t$,

$$\text{var}(G_t|N_t) = \bar{\Omega}_{t,2,2} - \frac{\bar{\Omega}_{t,1,2}^2}{\bar{\Omega}_{t,1,1}} \quad (62)$$