

Engineering Time Distributions in Gene Regulatory Networks

Seunghee Shelly Jang

I. INTRODUCTION

The role of stochasticity in biology has been studied in many contexts - such as creating variations required to survive in hostile environments [1], regulating circadian clocks [2], and probabilistic differentiation in developing cells [3]. If we systematically characterize the stochasticity in each context, we gain the ability to control these biological functions. The ability to control the development of multicellular systems, for one, holds fantastic futures like replacement organs grown *in vitro* or smarter drugs that only target tumors. With such potentials, then, which feature of stochasticity in development should we focus on?

In the development phase of multicellular organisms, an isogenic group of cells differentiates into multiple groups of heterogeneous cells with different epigenetics. This behavior can be likened to a group of people performing a leader election. To fairly elect a leader, the group can repeatedly and separately perform a task that has probabilistic outcomes, where one of the outcomes is the victory outcome. For example, a coin-toss with a head. After some number of tries, if an individual is the first one who ends up with a head among the group, he or she becomes a leader - effectively differentiating him or her from the group. The leader can then send signals to the rest of the group, telling them to stop flipping coins and become followers - again differentiating them from their undecided state, and from the new leader. Thus, it is not farfetched to imagine that a similar mechanism takes place in developmental differentiation.

Let us assume that an individual ends up with a head for the first time after h tries. Because coin-tosses have probabilistic outcomes, h is a random variable. And as such, h is characterized by its probability distribution. Let us assume that this probability distribution can be manipulated somehow - by biasing the coin, for example - then, the variance of h has an interesting interpretation in the leader election example. If the variance of h was set small, the probability of multiple undecided individuals each ending up with the victory outcome in a short amount of time near the average value is large - in other words, the victory outcomes are closely synchronized, and the group may end up with multiple leaders. However, if the variance is set large, the victory outcomes are asynchronous and the group is less likely to have multiple leaders. Therefore, the probability distribution of h affects the population distribution of differentiated states of leaders and followers.

In the following sections, we discuss the translation of the leader election example into a cellular context, specifically in the development phase, and propose a possible biological equivalent to the biased coin. Then we discuss the ways to characterize the probability distribution of h (or some equivalent random variable in the proper context), both with theoretical analysis by posing the scenario in mathematical language, and with experiments by synthesizing the biological biased-coin equivalent. This discussion of characterization methods is followed by the prerequisite backgrounds in both theory and experiments, as well as related works in the

field that serves as helpful starting points for the proposed research. Some preliminary results are discussed in the last section, along with recommendations and proposed future works.

Key words: Gene regulatory networks, synthetic biology, completion time, probability distribution, development

II. APPROACH AND OBJECTIVES

As long as the mechanism is probabilistic, one can suggest several intracellular environment analogs to the coin-flip example, such as protein dimerization, folding, or saturation. Let us consider the protein saturation example, where at $t = 0$, the gene coding for the protein of interest, X, gets activated and there is no X present. After activation, the gene begins to express and the count number of X (N_X) increases as long as the gene remains activated and the rate of X degradation is less than the rate of X synthesis. If the gene is deactivated before N_X reaches the saturation value, then N_X begins to decrease until the gene is activated again. This process of gene activation and deactivation occurs repeatedly until at some time $t = T_c$, the count number of X reaches the saturation value. This time T_c is analogous to h in the coin-flip example, and we call T_c the *completion time* of the protein synthesis process. And if we further assume that the count number of X determines the state of each cell, then the distribution of the differentiated states of the cell population with the gene X is affected by the probability distribution of T_c .

There can be a number of ways to manipulate the probability distribution of the completion time in this example. One way is varying the frequencies of gene activation and deactivation. Another way is varying the mechanism that activates the gene - an open-loop activation from external inputs or a feedback activation/deactivation by X. Frequency variations change the quantitative features (system parameter) of the gene regulation, and feedback or open-loop variations change the qualitative features (structure) of the gene regulation. The relationship between the quantitative and qualitative features of gene regulation and the probability distribution of the completion time will help us understand the fundamental design principles employed by nature to perform development and differentiation in multicellular organisms. Therefore, we propose the following objectives to guide the investigation of such relationship.

– **Synthesize single-gene networks in *E. coli*.** Three different mechanisms of gene regulation will be studied in this research - open-loop, positive feedback and negative feedback. The synthetic gene network corresponding to an open-loop mechanism will have a single gene that is activated by some external inputs. For the two feedback mechanisms, a single gene network that expresses either its own repressor or activator will be synthesized, to correspond to a negative or positive feedback mechanism, respectively. All of the gene networks will contain with inducible promoters and fluorescence protein gene. The inducible promoters allow us to measure the completion time by setting the initial time to when the promoter is induced, and the level of fluorescence

emitted by the fluorescence protein is measured to monitor the gene activity.

- **Mathematically model the three gene networks.** Using the Chemical Reaction Network theory, we will model the interactions among the gene network species. We apply a variety of stochastic analysis tools to the models in order to characterize the completion time, its probability distribution, and sensitivity to parameter variations and structural variations. Such analysis tools include the Chemical Master Equation (CME), the Stochastic Simulation Algorithm (SSA), and cumulant and moment dynamics. We will identify the qualitative differences of the gene networks arising from the difference in structure, and discuss how they can make each structure a better or worse suited mechanism used in development processes. In addition, quantitative features regarding the change in parameters will be investigated. The limitations on the probability distribution of completion time placed by physically feasible parameter values may also explain why a certain structure is more frequently observed in development than others.
- **Iteratively verify predictions made in the models with experiments and modify the models based on the experimental results.** The gene activity is monitored by measuring the level of fluorescence emitted by the synthesized fluorescence protein. The probability distributions of the completion time in these synthetic gene networks will be approximated using cellular assays, such as time-lapse microscopy or flow cytometry. Time-lapse microscopy allows us to monitor the individual trajectory of fluorescence level in a *single cell* and the time at which the fluorescence level reaches some saturation value. On the other hand, flow cytometry reveals the population *distribution* of fluorescence level at each measurement. Therefore, by measuring the fluorescence distributions at multiple times, we can study the distribution dynamics of the fluorescence level. And from the dynamics, we will derive the fraction of the population that has reached the saturation value at each measurement time, which leads to approximated probability distribution of computation time. The experimental results will be used to invalidate some of the candidate models and point out the features that require modification to attain better fidelity to the actual systems. The modified models, in turn, are used to design experiments that will better highlight the key features of the systems. The mathematical model predictions obtained from this iterative process will identify the salient features of differentiation mechanisms and allow us to synthesize gene networks with the complexity comparable to the naturally occurring examples.

The following section will provide a broad overview of the fundamentals in both theory and experiments to accomplish our objectives. Two specific related works are discussed afterwards, each with a focus on theory and experiments respectively. These works were chosen based on their close proximity to the objectives of the proposed research, and served as a foundation for obtaining the preliminary results that are discussed in a later section.

III. BACKGROUND AND LITERATURE REVIEW

As the biotechnology steadily advances, researchers are able to synthesize gene regulatory networks with increasing precision and success. These synthetic gene networks are built from borrowed biological components of natural genetic

regulatory parts, such as promoters and transcription factors. Though manipulating genetic materials is not a new technology, synthetic biology is different from the traditional genetic engineering in its intention to engineer novel behaviors, such as oscillation or bistability [4], [5]. The underlying objectives of these synthesis-based approach to biology is to identify and isolate the salient features of complex gene networks and discover the nature’s design principles. And synthetic biology is strengthened by two complementary approaches of mathematical theory and biological experiments. A well-established study of differential equations is used to analyze the dynamics of the systems [6], linear systems theory the stability and controllability [7], and probability theory the stochastic behaviors in the mesoscopic level of biological molecules [8], to name a few. At the same time, increasing efficiency of cloning techniques [9], decreasing cost of DNA synthesis and sequencing [10], and the advance of experimental equipments all contribute to engineering biological test beds for verifying hypotheses obtained from mathematical theories. As the objectives of the proposed research spans both theory and experiments, the rest of the background section is divided into two sections to address the fundamentals of each aspect separately.

A. Mathematical Theory

The theory of Chemical Reaction Network was originally developed to provide a standardized foundation from which a mathematical description of chemically interacting species inside a fixed volume can be derived [11]. The CRN of a given system contains chemical species (X_i) that interact with respect to some reaction (R_j), the stoichiometric coefficients of reactants (u_{ij}) and products (v_{ij}) of the chemical reactions, and the rates of these interactions (λ_j). From this description, using the Law of Mass Action, the dynamically changing concentrations of the chemical species are modeled by a set of ordinary differential equations. This method translates smoothly into the context of biological interactions inside a cell. Cellular environments are no different from the environments inside a chemical processing plant, such that they have biochemical interactions, reactant and product species of these interactions, and numerical values that describe the rates of the interactions. However, the key difference is that whereas chemical systems tend to have a large quantity of each species, the quantity of biological molecules tend to be present in much smaller quantities. Thus, chemical species can be expressed as continuous variables, but biological molecules must be expressed as discrete variables. Additionally, the stochasticity of the biochemical interactions become more pronounced in a system with species in small quantities. Therefore, biochemical systems, such as gene regulatory networks require mathematical description that properly addresses the *discrete* copy number of species and the *stochasticity* of interactions.

The discrete values of biochemical molecules inside gene regulatory networks allow us to model the systems as discrete-state continuous-time Markov processes [12]. Let the species of an arbitrary gene regulatory networks be denoted by a vector $\mathbf{S} = [S_1, \dots, S_n]$, and the number of each species are denoted by X_i . Each discrete state of the system is then denoted by the vector $\mathbf{X} = [X_1, \dots, X_n]$. Because the stochasticity of gene networks forces the description of the system from a singular deterministic value to a probability distribution over the states, we denote the probability of the system in state \mathbf{X} at time t by $p(\mathbf{X}, t)$. The vector of the probabilities of all the states is $\mathbf{p}(t)$ and the probability vector, given some initial

distribution, \mathbf{p}_0 , evolves according to the following master equation.

$$\dot{\mathbf{p}}(t) = \mathbf{Q}\mathbf{p}(t), \quad (1)$$

which is the Chemical Master Equation, and the matrix $\mathbf{Q} = [q_{ij}]$ contains the rates of system transitions from state j to state i [13]. The analytical solution of (1) is

$$\mathbf{p}(t) = e^{\mathbf{Q}t}\mathbf{p}_0. \quad (2)$$

Do not be misled by the elegantly simple form of the solution, as the matrix exponential, $e^{\mathbf{Q}t}$, requires an infinite sum of high computational cost. Instead of solving for the probability distribution dynamics analytically, a numerical alternative exists, where the exact realizations of the corresponding CRN can be obtained.

The Stochastic Simulation Algorithm (SSA) was developed to numerically simulate individual trajectories of the species of an arbitrary stochastic chemical reaction network [14]. The method employs the fact that 1) each rate of a chemical reaction is the inverse of the mean waiting time for the reaction, and 2) the probability of a reaction is equal to the ratio of the reaction propensity to the sum of all reaction propensities. Then, the time evolution of a stochastic system can be approximated by generating a large number of simulations of the system and studying the dynamics of each species. It should be clarified that the probability distribution dealt with in the CME corresponds to the *joint* probability distribution of each specific state, $P([X_1 = x_1, \dots, X_n = x_n], t)$, whereas the approximated probability distributions obtained via SSA are the *marginal* probability distributions of all species, $P(X_1, t)$. Though, more straightforward for portrayal of probability distribution dynamics, the SSA algorithm requires that the initial condition and the rate constants be specified *a priori*. It means that if an analysis requires a different set of rates or initial conditions, a whole new set of large number of simulations is required to study the specific condition.

A second alternative approach of characterizing the evolution of the probability distribution for stochastic biochemical systems is to compute the cumulant dynamic of each species of the system [15]. The cumulants of a random variable are set of values that characterizes the shape of the corresponding probability distribution. For example, the second order cumulant of a random variable is its variance and is representative of the width of the probability distribution. The cumulants are computed using the cumulant generator function,

$$G_X(s) = \log \langle e^{sX} \rangle, \quad (3)$$

where X is the random variable and $\langle \cdot \rangle$ denotes the expected value. The n th order cumulant of X is computed by taking the n th derivative of (3) with respect to s and setting $s = 0$. Usually, no more than the first four cumulants are computed for a given species, because cumulants of order five or higher have no straightforward interpretation related to the probability distribution characteristics. The time evolution of these cumulants requires an additional function called the extended generator. Let $\psi(X(t))$ be some test function of state $X(t)$, then the expected value of this test function evolves according to the following equation.

$$\begin{aligned} \frac{d \langle \psi(X(t)) \rangle}{dt} &= \langle L\psi(X(t)) \rangle \\ &= \sum_{j=1}^m \lambda_j (\psi(X^j(t)) - \psi(X(t))), \quad (4) \end{aligned}$$

where the $X^j(t)$ is the state after the reaction $\mathbb{R}_j: X(t) \mapsto X^j(t)$ has occurred, λ_j is the reaction rate constant, and L is the extended generator. The cumulant dynamics is then obtained by letting $\psi(X(t)) = G_X(s)$, and solving the resulting set of ordinary differential equations. An interesting connection exists between the cumulant dynamics and the differential equation obtained by using the Law of Mass Action, such that the first order cumulant dynamics is equal to the deterministic dynamics predicted by Mass Action kinetics. In fact, this is not surprising because the first order cumulant is the mean of the population and Law of Mass Action predicts the average behavior of the population.

B. Biological Backgrounds

Within a single cell resides a genome, a chain of DNA molecules, that contains all the genetic information the cell needs to harvest energy, reproduce and survive. The genome alone, though mighty in its information content, cannot make a living organism. It requires molecular machinery that actualizes this information in useful form, thus is the function of RNA and protein. DNA is transcribed into RNA, and in turn the RNA is translated into protein, and proteins are the true workers of biological functions [16]. The most critical role of protein molecules is regulating genomic information transfer as it affects the viability of an organism. If the processes of transcription and translation were not properly regulated, in other words if the entire genome was uniformly transcribed and translated, it would mean a disaster for the cell. Therefore, there exists intricately connected networks of gene regulation that allows cells to allocate energy, respond to its environment and procreate.

The two major components of gene regulatory mechanisms are promoters and transcription factors (TF). TFs are protein complexes that act either as a repressor or an activator by binding to the promoter of a gene. Promoters are short sequence of DNA that are located at the 5'-end of a gene and are recognized by RNA polymerase to initiate an RNA synthesis. A bacterial promoter has two short 6 basepair long sequences that are conserved in most promoters, called the consensus sequences. The rest of the promoter sequences are composed of operators that serve as binding sites for specific TFs. A large number of TF and promoter pairs have been identified in metabolic pathways, signal transduction pathways, and developmental regulatory pathways. The known pairs of TF and promoter are used to design and build synthetic gene regulatory networks by arranging them in specific configurations [17], [18]. For example, the critical structure of stress response in *B. subtilis* were identified by synthesizing a gene network with identical promoters and TFs, but with one of the two feedback loops (coupled positive and negative feedbacks) removed [19]. This synthetic version, when transformed inside cells, prohibited the cells from leaving their competence state, thus showing that the removed feedback is critical to the overall mechanism of *B. subtilis* stress response.

Aside from this example, there are plenty of examples where feedback mechanisms are observed. A class of gene networks that give rise to stochastic state switching, such as cancer and developmental differentiation, has been consistently shown to contain positive feedback loops [20], [21], [22], [23]. Another type of behavior that arises from gene networks with positive feedback loops is procrastinating differentiation [24]. Procrastination refers to the phenomenon observed in isogenic cells, that when triggered for specific response (e.g.

sporulation, apoptosis), the response time of each cell widely vary within the microcolony. This phenomenon is closely related to the coin-flip leader election example - however, the lack of communication between the individuals presents itself in the form of non-differentiated steady-state, i.e. the population distribution is still uniform, and the variance is only observed during the transient process.

C. Biological Experiments

Gene network synthesis procedure can be broken into two major steps. The first step is acquiring the desired DNA sequences (e.g. promoters, transcription factor genes), and the second step is joining these pieces together in the right order. Natural promoters and genes are obtained from the host organism's genome through Polymerase Chain Reaction (PCR). In this process, the desired sequence is isolated and amplified by using two short pieces of single-stranded DNA that are complementary to the 3'- and 5'-ends of the desired sequences, and DNA polymerase. Through cyclical temperature manipulation of the reaction chamber, the quantity of the desired sequence amplifies exponentially. Then the amplified pieces are digested using restriction enzyme to introduce recognizable sticky (or blunt) ends to each piece. The individual pieces are assembled together using DNA ligases, which are originally derived from DNA repair mechanism. The restriction enzymes were chosen strategically to ensure that when the DNA ligase is in action to assemble the pieces together, they are joined in the right order and direction. The product is transformed into a host organism by electroporation in the last step before it is ready for assays.

The first usage of green fluorescent protein (GFP) as a marker for gene expression significantly changed the way gene activity is monitored [25]. Modern cellular assay tools such as the flow cytometry, time-lapse fluorescence microscopy, and microplate readers are made possible by GFP-fusion proteins. GFP is fused to a protein of interest via DNA recombination and is monitored for fluorescence emission upon excitation. GFP molecules have excitation wavelength of 395nm or 498nm, depending on the organism it was derived from, and emission wavelength of 505nm. In addition to GFP, there are several different types of fluorescence protein that can be used for multicolor live cell imaging to monitor several different gene activities simultaneously.

D. Related Works

One way of deriving the analytical expression of completion time probability distribution is to solve the CME of the system in Laplace domain [26]. In this work, a kinetic proofreading (KPR) process was modeled by a Markov chain with an absorbing state, where the absorbing state corresponded to the completion of the proofreading process that required sequential intermediate steps. In recognition that the completion time is essentially the first-passage time of Markov chain, they performed Laplace transform to the solution of the CME shown in (2) to obtain the analytical expression [27]. The solution showed that the distribution of the first-passage time approaches limiting behaviors, depending on the direction of the bias imposed by the transition rates - forward to the sink state, or backward to the initial state. However, the solution and the conclusion is limited to an open-loop system where the transition rates are independent of the states. Though the authors analyze simulated systems with *varying* transition rates by using randomly generated values and expand their conclusion, this is quite different from *state-dependent* transition rates of feedback mechanisms. It will be

interesting to investigate whether a similar conclusion can be drawn from biochemical processes with feedback.

A feedback loop in a gene regulatory network consists of a promoter that are regulated by some TF, which in turn, is expressed by the gene controlled by the promoter. Naturally occurring feedback loops are interesting in themselves, but a synthetic class of hybrid promoters developed to exhibit the programmability of promoters expands the number of possible feedback loops [28]. Hybrid promoters are synthesized by combining multiple operator sites corresponding to different TFs, so that the resulting promoters are regulated by more than one type of TF. Additionally, the order in which these operator sites are arranged was shown to affect the expression level of the downstream gene. All the hybrid promoters studied in [28] are inducible by specific inducer chemicals and the concentrations of these inducers are shown to be correlated with the downstream gene expression as well. Therefore, by employing these hybrid promoters and varying the inducer concentrations and the copy number of genes, many variations of feedback loops can be engineered.

IV. PRELIMINARY RESULTS

A. Synthetic positive feedback gene network in *E. coli*

Design and Construction. Since even a small amount of output can trigger the rapid activation of a positive feedback mechanism, we employed promoters that have minimal leaky expression and give maximal control over the range of expression [29]. We selected two hybrid promoters from [28] - A12 and D61 - that have operator sites from the pBAD promoter that is activated by AraC-arabinose complex, and the pLac promoter that is repressed by LacI protein. To create a positive feedback loop, we cloned *araC* gene downstream of the promoter so that when induced with arabinose, the promoter is activated (Figure 1 (a)). The cloning procedure was conducted using restriction enzymes and DNA ligase. Additionally, we used a strain of *E. coli* that constitutively express *lacI* to transform the construct into, so that the repression of *LacI* keeps the promoter tightly regulated when uninduced. Currently, there are several variations of the positive feedback gene network (Figure 1 (b)). Each network has the promoter A12 or D61, and placed in a plasmid with the origin of replication pMB1, pSC101, or pSB3K3¹.

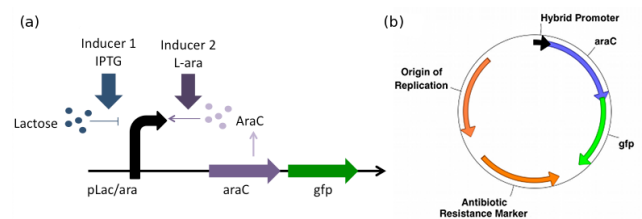


Fig. 1: A diagram of a synthetic positive feedback gene regulatory network. (a) The pLac/ara hybrid promoter is induced by IPTG and arabinose. IPTG inhibits the *LacI* repression of the promoter, whereas arabinose forms a complex with *AraC* protein and activates the promoter. (b) The network has a hybrid promoter, controlling the expression of *araC* and *gfp* downstream, an origin of replication and an antibiotic resistance marker.

Assays. Three different types of assay were used to characterize the A12 promoter variant of the positive feedback gene networks in varying concentrations of IPTG and arabinose. The following list provides a brief summary of each assay, including objectives, method, results, and conclusion.

¹Each has an approximate copy number of 10^0 , 10^1 , and 10^2 , respectively

– plate spectrophotometry

- *Objective* - To confirm the response behavior of the hybrid promoter predicted in the original paper in changing concentrations of two inducers.
- *Methods* - We chose a single colony from the agar plate of transformed cells and inoculated LB media. The culture was diluted 1:150 in PBS with 48 different concentrations of IPTG and arabinose in a 96 well plate (duplicates were made for each condition). Each well contained 0%, 0.01%, 0.05%, 0.1%, 0.5%, 1% or 2% arabinose and 0uM, 10uM, 50uM, 100uM, 500uM, or 1mM IPTG. The plate reader was set to measure the optical density and the fluorescence level of each well every 20 minutes and shake for 24 hours. Additionally, the plate was kept in 37C and shaken for 10 minutes before each measurement.
- *Results and Discussion* - The A12 variant of positive feedback gene network had minimal change in the steady-state fluorescence value with varying concentration of arabinose, but significant change with varying concentrations of IPTG. This confirmed that the hybrid promoter response was consistent with the prediction.

– flow cytometry

- *Objective* - To measure the population distribution of fluorescence level of cells from the time they are induced until steady-state is reached, and to observe the transient in the mean and the variance of the fluorescence level.
- *Methods* - We used the same colony to inoculate LB cultures with varying concentrations of IPTG and arabinose. 12 different concentrations were tested, and each culture tube contained 0%, 0.1% 1% or 10% arabinose and 0uM, 10uM or 100uM IPTG. In 15 minute intervals, for 4.5 hours, 10uL of culture from each tube was diluted 1:15 in 96 well plates with 140uL of PBS in each well. 50000 events (cells) from each well were screened using a Accuri C6 flow cytometer. Additional measurements was made at 5 hrs, 6 hrs, 19 hrs and 20.5 hrs after induction.
- *Results and Discussion* - At 30 minutes after induction, the mean of the fluorescence level distribution was higher compared to the initial mean value at the time of induction. From 45 min to 105 min after induction, the mean fluorescence level decreased. At 120 min after the induction two distinct populations of fluorescence level were observed and again at 120 min. However, the bimodal distribution disappeared abruptly in the next measurement and no more higher fluorescence population was observed. We hypothesize that the cell population with activated positive feedback are suffering from AraC/gfp toxicity.

– time-lapse microscopy

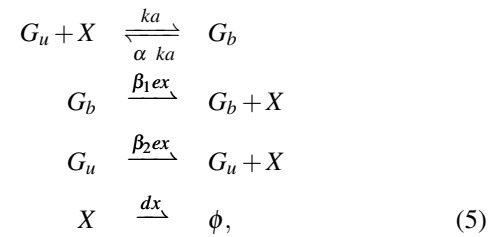
- *Objective* - To confirm whether the high fluorescing cells - the cells with activated positive feedback - have different viability compared to the lower fluorescing cells.
- *Methods* - An agar plate with 10% arabinose and 1mM IPTG was prepared. An overnight culture was diluted in the morning and grown for 3 hours to reach log-growth phase. xxxuL of cells were transferred onto the agar plate. Using the microscope, 10 sparsely populated areas were selected. Within each area, a single cell was marked for tracking. Every 10 minutes an image processing macro tracked each cell, adjusted the focus, and took an image of the cell, over 12 hours.
- *Results and Discussion* - The cells that began to emit high level fluorescence were shown to grow larger than its peers with low level fluorescence. Eventually these cells died,

adding to the suspicion of AraC/gfp toxicity resulting from the amplification effect of positive feedback loops.

To confirm whether the cells are affected by araC toxicity, we will transform the gene network into strain of *E. coli* that has its araC gene from its genome removed [30]. Then similar assays will be conducted to find the optimal concentrations that will delay the accumulation of AraC. Additionally, a different hybrid promoter with lower level of expression will be switched in to test whether a better viability of the cells can be obtained.

B. Approximation of the probability distribution of completion time

Using the basic understanding of gene expression mechanism, we expressed the feedback gene regulatory network with the following Chemical Reaction Network.



where G_u is an X gene not bound with the transcription factor X, G_b is an X gene bound with X, and X is the transcription factor. The rate constants of reactions are; rate of transcription factor binding (ka), unbinding-to-binding ratio (α), ratio of unbound gene expression to unregulated expression (β_1), ratio of bound gene expression to unregulated expression (β_2), and transcription factor degradation/dilution rate (dx). We simulated 1000 SSA realizations of the CRN to visualize the dynamic of X and approximate the probability distribution of the completion time in this specific condition (Figure 2). We set the completion of protein saturation as when the X reaches its half steady-state value (N), because it was observed that after an individual simulation reaches X value of N, the probability of the X value dropping below N is small. In fact, we have not seen such event in our limited number of simulation, and it is a generous over-estimation. We can be convinced with some confidence that the cell is committed to its fate, after the count number exceeds N. However, for future work, it is recommended that by modeling the system as Markov processes and rigorously determine the probability of the system returning to its initial state after reaching the state of $X \geq N$.

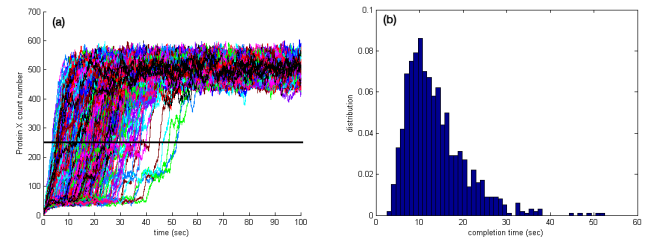


Fig. 2: (a) Thousand trajectories of the CRN obtained from the Stochastic Simulation Algorithm. (b) the approximated probability distribution of completion time (completion is when $X > N$, the horizontal black line in (a)). The parametric values used are $[ka, \alpha, \beta_1, \beta_2, ex, dx] = [10^{-3} \log(2), 0.1, 10, 1, 10 \log(2), \log(2)]$ and the initial condition is $[G_{unbound}, G_{bound}, X] = [5, 0, 0]$.

We applied the extended generator method to compute the cumulant dynamics of the system. However, the CRN involves reactions of order higher than 1 (i.e. more than one reactant in

a given reaction). When the differential equations of species cumulants are derived, cumulants of order i is a function of the next order cumulants, $i + 1$. Therefore, to obtain the exact solution, an infinite order of cumulants are required. Therefore, we truncate the cumulants with order 3 and higher, in order to get a closed solution of the set of differential equations. This is equivalent to assuming that the population of each species has a Gaussian distribution at all times. A Gaussian distribution is a well defined function,

$$g(X) = \frac{1}{\sqrt{2\pi\kappa_{XX}}} \exp - \frac{(X - \kappa_X)^2}{\sqrt{2\kappa_{XX}}}, \quad (6)$$

where κ_X and κ_{XX} are the first and the second order cumulants of X . The function allows us to approximate the fraction of population located above the saturation value, N , as a function of time, using the erf function

$$F(N,t) = \frac{1}{2} \pm \frac{1}{2} \operatorname{erf} \left(\frac{N - \kappa_X(t)}{\sqrt{2\kappa_{XX}(t)}} \right) \quad (7)$$

This function turns out to be an approximation of the cumulative distribution of the completion time. By taking the derivative of this function, we approximate the probability distribution of completion time.

$$\begin{aligned} f(N,t) &= \frac{d}{dt} F(N,t) \\ &= \frac{1}{2} \left(\frac{-\sqrt{\kappa_{XX}(t)} \kappa'_X(t) - \kappa'_{XX}(t)(N - \kappa_X(t))}{\kappa_X(t)} \right) \times \\ &\quad \exp \left(- \frac{(N - \kappa_X(t))^2}{\kappa_{XX}(t)} \right) \end{aligned} \quad (8)$$

Though a fair approximation, the function lacks certain mathematical rigor, because the normal distribution assumption only applies to continuous variables. A similar approach where the first and second order cumulants are sufficient to approximate the entire *discrete* variables is required.

V. PLAN OF WORK

A. Experiments

The concentration variation assays discussed in the Preliminary results section is currently incomplete for the other hybrid promote, D61. Considering the suspected AraC/GFP toxicity resulting from positive feedback in the A12 hybrid promoter gene network, we expect a similar result in cellular assay - death of cells with high fluorescence level. However, since D61 has lower expression level compared to A12 (about an order of magnitude lower), it is possible that the cells with activated positive promoter will suffer less.

We plan on running a series of single-cell tracking time-lapse fluorescence microscopy of the gene networks in varied concentrations of inducer to obtain time series data of fluorescence level that resembles Figure 2 (a), and subsequently an approximated completion time probability distribution that resembles Figure 2 (b). The mean and the variance of this distribution are expected to vary with changing conditions of the inducers. And changing this value is closely related to varying the parameters ka and α in the model (5). Additionally, we will clone a number of other hybrid promoters in the current gene network, to see the effect of β_1 and β_2 variations in the completion time probability distribution.

In addition to the positive feedback gene networks, we will design and analyze negative feedback gene networks and open-loop gene networks. The focus of the design is to keep the details, such as the inducers, the plasmids, and the

E. coli strains, as equivalent as possible, and vary only the structure of the feedback. This is so that there are no external bias affecting the assay results. The analysis process will be consistent with that of the positive feedback gene networks (i.e. plate spectrophotometry, flow cytometry, and time-lapse microscopy).

B. Theory

The analytical solution of the completion time probability distribution will be explored using a variety of tools (matrix norms, asymptotics, etc). A better way of approaching the problem might be in identifying the upper and lower bounds of the variance and mean of the completion time, instead of identifying the exact solutions. Once the analytical solutions are obtained for the bounds, the parametric sensitivity of these values will be investigated. The purpose of this analysis is to identify the parameters that would result in the largest observable difference in the experiments with the synthetic gene regulatory system. Additionally, we will study what possible type of completion time probability distributions can be engineered using gene regulatory networks (e.g. Gaussian, uniform, etc).

As with the experimental objectives, we will analyze negative feedback and open-loop systems for the same features, completion time distribution. We expect to observe qualitative difference in the shape of the probability distribution of completion times in these systems, and will discuss the reasons for these differences. Furthermore, we will derive the connection between the frequency of observation of each type of feedback in development/differentiation process and the qualitative difference of the completion time probability distribution.

VI. SCHEDULE AND REQUIRED RESOURCES

The schedule of work is shown in Figure 3. The biological equipments and computational software required for research are funded by the Molecular Programming Project, part of the National Science Foundation's Expedition in Computing program.

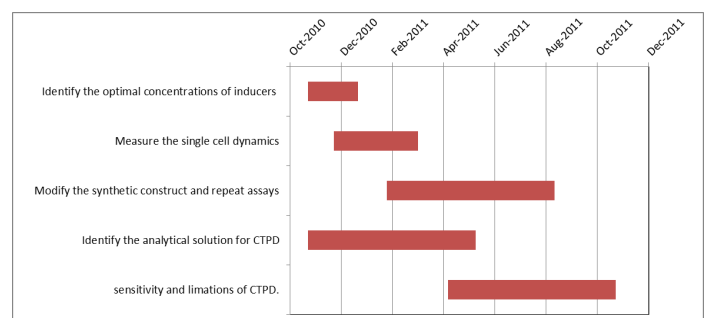


Fig. 3: Schedule of work

REFERENCES

- [1] D. Wolf, V. Vazirani, and A. Arkin. Diversity in times of adversity: probabilistic strategies in microbial survival games. *Journal of Theoretical Biology*, 2005.
- [2] M. Kaern, T. Elston, W. Blake, and J. Collins. Stochasticity in gene expression: from theories to phenotypes. *Nature Reviews Genetics*, 2005.
- [3] T. Suda, J. Suda, and M. Ogawa. Single-cell origin of mouse hemopoietic colonies expressing multiple lineages in variable combinations. *PNAS USA*, 80:6689, 1983.
- [4] M. Elowitz and S. Leibler. A synthetic oscillatory network of transcriptional regulators. *Nature*, 403, 2000.
- [5] T. Gardner, C. Cantor, and J. Collins. Construction of a genetic toggle switch in *Escherichia coli*. *Nature*, 403(20), 339-342 2000.
- [6] M. Hirsch, S. Smale, and R. Devaney. *Differential equations, Dynamical systems, and an introduction to chaos*. Elsevier Academic Press, 2004.

- [7] N. Nise. *Control Systems Engineering*. Wiley, 2000.
- [8] G. Grimmett and D. Stirzaker. *Probability and Random Processes*. Oxford University Press, 2001.
- [9] D. Gibson and et. al. Creation of a bacterial cell controlled by a chemically synthesized genome. *Science*, 2010.
- [10] R. Carlson. The changing economics of DNA synthesis. *Nature Biotechnology*, 27, 2009.
- [11] M. Feinberg. Lectures on Chemical Reaction Networks. Online - Departments of Chemical Engineering and Mathematics, The Ohio State University, 1979.
- [12] D. Stroock. *An Introduction to Markov Processes*. Springer, 2005.
- [13] D. Gillespie. A rigorous derivation of the chemical master equation. *Physica A: Statistical Mechanics and its Applications*, 188:404–425, 1992.
- [14] D. Gillespie. Exact Stochastic Simulation of Coupled Chemical Reactions. *The Journal of Physical Chemistry*, 81(25), 1977.
- [15] J. Hespanha. Modeling and Analysis of Stochastic Hybrid Systems. *IEEE Proc - Control Theory and Applications, Special Issue on Hybrid Systems*, 153:520–535, 2007.
- [16] B. Alberts. *Molecular Biology of the Cell (Chapters 1 - 7)*. New York : Garland Science, 4th edition, 2002.
- [17] M. Elowitz, A. Levine, E. Siggia, and P. Swain. Stochastic gene expression in a single cell. *Science*, 297:1183–1186, 2002.
- [18] T. Knight. Idempotent Vector Design for Standard Assembly of Biobricks. Technical report, MIT, 2003.
- [19] G. Suel, J. Garcia-Ojalvo, L. Liberman, and M. Elowitz. An excitable gene regulatory circuit induces transient cellular differentiation. *Nature*, 440:545–550, 2006.
- [20] A. Kashiwagi, I. Urabe, K. Kaneko, and T. Yomo. Adaptive Response of a Gene Network to Environmental Changes by Fitness-Induced Attractor Selection. *PLoS ONE*, 2006.
- [21] M. Levine and E. H. Davidson. Gene regulatory networks for development. *Proc. Natl. Acad. Sci. U.S.A.*, 102:4936–4942, Apr 2005.
- [22] E. H. Davidson, J. P. Rast, P. Oliveri, A. Ransick, C. Caletani, C. H. Yuh, T. Minokawa, G. Amore, V. Hinman, C. Arenas-Mena, O. Otim, C. T. Brown, C. B. Livi, P. Y. Lee, R. Revilla, A. G. Rust, Z. Pan, M. J. Schilstra, P. J. Clarke, M. I. Arnone, L. Rowen, R. A. Cameron, D. R. McClay, L. Hood, H. Bolouri, and E. H. Davidson. A genomic regulatory network for development. *Science*, 295:1669–1678, Mar 2002.
- [23] V. F. Hinman, A. T. Nguyen, R. A. Cameron, E. H. Davidson, and E. H. Davidson. Developmental gene regulatory network architecture across 500 million years of echinoderm evolution. *Proc. Natl. Acad. Sci. U.S.A.*, 100:13356–13361, Nov 2003.
- [24] A. Eldar and M. Elowitz. Functional roles for noise in genetic circuit. *Nature*, 467, 2010.
- [25] M. Chalfie, Y. Tu, G. Euskirchen, W. Ward, and D. Prasher. Green Fluorescent Protein as a Marker for Gene Expression. *Science*, 263(5148):802–805, February 1994.
- [26] G. Bel, B. Munsky, and I. Nemenman. The simplicity of completion time distributions for common complex biochemical processes. *Physical Biology*, 7(1), 2010. Bel, Golan Munsky, Brian Nemenman, Ilya.
- [27] S. Redner. *A Guide to First-Passage Processes*. Cambridge University Press, 2001.
- [28] R. Cox, M. Surette, and M. Elowitz. Programming gene expression with combinatorial promoters. *Molecular Systems Biology*, 2007.
- [29] R. Lutz and H. Bujard. Independent and tight regulation of transcriptional units in *Escherichia coli* via the LacR/O, the TetR/O and AraC/I1-12 regulatory elements. *Nucleic Acids Research*, 1997.
- [30] T. Baba, T. Ara, M. Hasegawa, Y. Takai, Y. Okumura, M. Baba, K. A. Datsenko, M. Tomita, B. L. Wanner, and H. Mori. Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol. Syst. Biol.*, 2:2006.0008, 2006.