

Gene regulatory network drift and speciation occurs rapidly under neutrality

Joshua S. Schiffman[†] Peter L. Ralph^{†‡}

[†]University of Southern California, Los Angeles, California [‡]University of Oregon,
Eugene, Oregon
jsschiff@usc.edu plr@uoregon.edu

Abstract

Here we introduce an analytical theory to study the evolution of biological systems, such as gene regulatory networks. The evolutionary conservation of phenotype under selective and environmental stasis does not necessitate conservation of the underlying mechanism, as distinct molecular pathways can realize identical phenotypes. Here we give an exact expression for the set of all linear mechanisms with identical phenotypes, and expect evolution under neutrality to explore this set. We employ a quantitative genetic approach to model evolution under neutrality as a random process over the set of all phenotype-invariant mechanisms: only mutational tweaks to the pathway that leave the phenotype invariant are optimally fit. We show that there is never a unique linear system architecture for any phenotype and that the evolutionary exploration of these distinct and mutationally connected mechanisms can lead to the rapid accumulation of hybrid incompatibilities between allopatric populations and thus lead to the rapid formation of new species – in fewer generations than there are breeding individuals in a population.

Additional ideas to consider adding:

- hybrid vigor (need to do calculation)
- discuss linearity, linearization, and canalization in introduction

Note: in the L^AT_EX source I'm putting in semantic linebreaks, so it's easy to edit and move around phrases and ideas.

Need to come up with a consistent term for “the A_{ij} ”s. – “regulatory coefficients”? “genotype”?

Introduction

Bridging the gulf between an organism's genome and phenotype is a poorly understood and complex molecular machinery. Progress in a suite of biological subdisciplines is stalled by our general lack of understanding of this molecular machinery: with respect to both its function and evolution. There does exist a growing body of data on the evolutionary histories and molecular characterizations of particular gene regulatory networks [??], as well as thoughtful verbal and conceptual models [????]. However, as Hardy and Weinberg taught us over a century ago, verbal theories are often insufficient, if not downright misleading [???]. This is especially pertinent given the staggering complexity and scope of contemporary research programs. This outlook necessitates the advancement of conceptual frameworks of such precision, only mathematics will suffice, as models allow the development of concrete numerical predictions. Previously it has been suggested that any idealized study of evolution is incomplete without a mathematically sufficient description of the genotype, phenotype, and transformation from one to the other [?].

Saying “stalled” reflects negatively on other fields - rephrase?

I don't think the HW example will speak to our audience. Also, better to say why math is good rather than why not-math is bad.

The molecular machinery, interacting with the environment, and bridging genotype to phenotype can be mathematically described as a dynamical system – or a system of differential equations [?]. Movement in this direction is ongoing, as researchers have begun to study the evolution of both abstract [????] and empirically inspired computational and

probably some more recent Wagner papers in this line as well?

mathematical models of gene regulatory networks (GRNs) [????????]. If we allow the reasonable assumption that the genotype-phenotype map can be represented as a system of differential equations, we can immediately discuss its evolution and function in a much more mechanistic, yet general, manner.

In some fields that seek to fit parametric models to experimental data, such as control theory, chemical engineering, and statistics, it is well known that mathematical models can fundamentally be *unidentifiable* and/or *indistinguishable* – meaning that there can be uncertainty about an inferred model’s parameters or even its claims about causal structure, even with access to complete and perfect data [???]. Models with different parameter schemes, or even different mechanics can be equally accurate, but still not *actually* agree with what is being modelled. In control theory, where electrical circuits and mechanical systems are often the focus, it is understood that there can be an infinite number of “realizations”, or ways to reverse engineer the dynamics of a black box, even if all possible input and output experiments on the black box are performed [???]. In chemical engineering, those who study chemical reaction networks sometimes refer to the fundamental unidentifiability of these networks as “the fundamental dogma of chemical kinetics” [?]. In computer science, this is framed as the relationship among processes that simulate one another [?]. Although this may frustrate the occasional engineer or scientist, viewed from another angle, the concepts of unidentifiability and indistinguishability can provide a starting point for thinking about externally equivalent systems – systems that evolution can explore, so long as the parameters and structures can be realized biologically. In fact, evolutionary biologists who study convergent versus parallel evolution, homology, and analogy are very familiar with such functional symmetries; macroscopically identical phenotypes in even very closely related species can in fact be divergent at the molecular and sequence level [????????].

In this paper we outline a theoretical framework to study the evolution of biological systems. We expect wide applicability to this approach (*i.e.* non-neutral evolution), however presently, we focus solely on neutral evolution, that is where phenotype is conserved over evolutionary time.

We derive an analytical description of the set of all linear biological systems with identical phenotypes – that is we describe the set of all gene network architectures that yield identical phenotypes, and show that all biological systems can, in principal, can undergo systems drift. In the neutral case, this set describes a manifold that evolution explores leaving phenotype invariant with respect to mutation, and predicts that if two populations become reproductively isolated, hybrid incompatibility can occur, despite the absence of adaptation, directional selection, or environmental change. Speciation typically occurs on timescales approximately on the order of N_e generations, where N_e is the effective population size.

Gene Networks as Linear Dynamical Systems

Here we outline a method to model biological systems, such as gene regulatory networks, as linear dynamical systems. We define an organism’s “phenotype” to be the time dynamics of molecular concentrations directly relevant to survival, and denote by $\phi_k(t)$ the concentration of the k^{th} component of phenotype at time t . These dynamics are the result of the interconnections of a gene regulatory network (or some other biological system) A (an $n \times n$ matrix) and an environmental input $u(t)$. Such a system \mathcal{S} , is composed of two equations,

$$\mathcal{S} := \begin{cases} \dot{\kappa}(t) &= A\kappa(t) + Bu(t) \\ \phi(t) &= C\kappa(t) \end{cases} \quad (1)$$

Where the *kryptotype* $\kappa(t)$, is a list of a system’s molecular concentrations at time t – or more generally the systems *internal dynamics*. As some system dynamics may not be visible

not clear what you mean by “actually”

here you are jumping to evolution – haven’t said yet why evolution “can explore” these – needs to say explicitly what we mean by neutral here

rather than “we expect” maybe say that the framework could be applied to non-neutral evolution

here’s the definition of neutral, but it’s not quite right?

We don’t “derive” it, we point out that Kalman already derived it.

We need something here saying what we mean by “gene network architecture”. Also, I don’t think we claim anything for “all biological systems”.

Below we mix general language with language specific to GRNs, like switching between “internal state” and “transcription factor concentration. Rather than keep making comments like “(or some other system)”, we could say up top that this applies to other situations, but to make it concrete we’ll talk about regulatory networks.

tried to say it more precisely. better phrase than “component of phenotype”?

say what u is (a vector)

say in words what κ is before the equation.

nor relevant to selection, we distinguish the *kryptotype* (as it is hidden) from the *phenotype* – which is the systems *external dynamics*. The change in internal molecular concentrations as a function of time $\dot{\kappa}(t)$ is simply a function of the systems current state (or kryptotype), the system architecture A – such as a gene regulatory network, the environmental input $u(t)$, and B (an $n \times l$ matrix) – how the system processes its environment. C (an $l \times n$ matrix) filters the dynamics relevant to survival. Thus the phenotype is simply a convolution of the system organization and the environment,

$$\phi(t) = \int_0^t C e^{At} B u(t) dt, \quad (2)$$

where we sometimes refer to $h(t) := C e^{At} B$ as the *impulse response* of the system.

As A specifies gene regulatory network architecture, the i^{th} row of A can be interpreted as the *cis*-regulatory element (or promoter) for gene i .

Example 1 (Oscillating Gene Network: Cell Cycle Control) *Cellular division is a complicated phenomena, governed by many different processes, however it is agreed that its rhythm is partially controlled by periodic (oscillating) gene transcription [?]. Consider a simplified model of oscillating gene transcription. In the present framework periodic expression requires at minimum two interacting genes.*

Suppose gene-2 up-regulates the transcription of gene-1 and that gene-1 down-regulates the transcription of gene-2 with equal magnitude (of 1) and relative to each of their concentrations, denoted by κ_1 and κ_2 . Furthermore, suppose that only the dynamics of gene-1 are consequential to the cell cycle (perhaps the amount of gene-1 activates another downstream gene network). Lastly suppose that the production of both genes is stimulated by an impulse of a molecule present immediately after division.

If the rate each of these genes is expressed is a linear function of their concentrations, the dynamics of the system are given by

$$\begin{aligned} \dot{\kappa}_1(t) &= \kappa_2(t) + u(t) \\ \dot{\kappa}_2(t) &= -\kappa_1(t) + u(t) \end{aligned}$$

where $\dot{\kappa}$ denotes the time derivative. The initial conditions $\kappa_1(0)$ and $\kappa_2(0)$, and the input $u(t)$ then determine the concentrations through time. If we record the regulatory coefficients in the matrix

$$A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix},$$

and define the column vector $B = [1 \ 1]^T$, then in matrix notation the dynamics are

$$\dot{\kappa} = A\kappa(t) + Bu(t).$$

Since only the dynamics of gene-1 are directly relevant to biological function, the dynamics of interest are given by

$$\phi(t) = C\kappa(t)$$

where the row vector C is defined as $C = [1 \ 0]$. (If the dynamics of both genes were physiologically relevant to the cell cycle, we would set C to be the identity matrix).

Since the input is simply an impulse, its phenotype is equivalent to its impulse response

$$\phi(t) = h(t) = \sin(t) + \cos(t).$$

Could say instead of "filters" that e.g. C_{ij} is the amount by which the j th input u_j affects the production of the i th transcription factor concentration.

This "Thus" needs expanding, e.g., "the solution to this equation is unique and given by" with a reference.

You've missed out $\kappa(0)$. Also, you can't have t be both a limit of the integral and the variable being integrated over – check this and get it right.

How about expanding this minor comment into a short paragraph with more interpretation of not just A but also B and C ? Should say why. Include talking through why $A_{ij} > 0$ means that j "upregulates" i .

you haven't explained $u(t)$ above. and, isn't it more natural in this set-up to have $B = [1, 0]$ than $[1, 1]$?

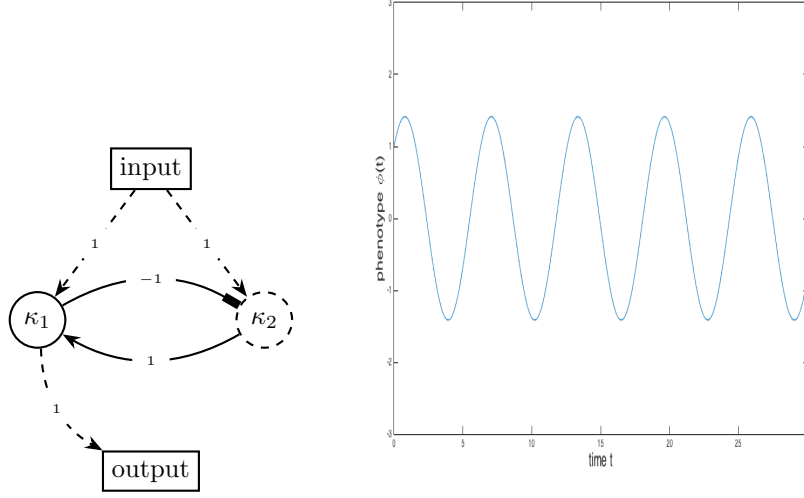


Figure 1: (Left) Graphical representation of the cell cycle control gene network, and (right) plot of the phenotype $\phi(t)$ against time t .

We return to the evolution of such a system below.

Linear Evolutionary Systems

The literature is filled with detailed observations of molecular systems and their diversity. There are examples of significant diversity in the networks underlying processes such as circadian rhythm [?], cell cycle control [??], pattern formation, and metabolism [?????]. Despite a symmetry in functionality or phenotype these systems often differ, sometimes substantially, at the molecular level. How many different mechanisms have the same function? We urge the reader to consider the metaphorical black box.

Gene regulatory networks with identical phenotypes (external dynamics) do not necessarily have identical kryptotypes (internal dynamics). Any linear and minimal system (a gene network) – minimal, informally meaning that the system’s phenotype is achieved with the fewest possible number of genes – has an identical phenotype up to a change of coordinates.

$$h(t) = Ce^{At}B \quad (3)$$

$$= CV^{-1}Ve^{At}V^{-1}VB \quad (4)$$

$$= CV^{-1}e^{VAV^{-1}t}VB \quad (5)$$

$$= \bar{C}e^{\bar{A}t}\bar{B} \quad (6)$$

Two biological systems, $\mathcal{S} = \{A, B, C\}$, and $\bar{\mathcal{S}} = \{\bar{A} = VAV^{-1}, \bar{B} = VB, \bar{C} = CV^{-1}\}$, have the same phenotype if they are related by a change of coordinates.

Although systems may not be identifiable beyond a change of coordinates, at present we are primarily interested in a subset of these systems. That is, systems that not only have equivalent external dynamics, but also equivalent input and output relationships. Formally, this means systems related by a change of coordinates (any invertible matrix V) that leaves

instead title this section "Systems drift for linear systems"?

This paragraph feels like the introduction.

Diversity doesn't imply systems drift – only if the diverse systems are homologous.

cite?

"has an identical phenotype" is not what you mean to say here

what are we calling these – not "biological system" – that's too broad

is this the place to say "and only if, in the minimal dimension"?

clarify that you mean actually having the same B and C . (at first I read "input and output relationship" as the mapping $u \rightarrow \phi$), i.e. behaviour across many possible inputs.

B and C invariant:

$$VB = B \implies \bar{B} = B \quad (7)$$

$$CV = C \implies \bar{C} = C \quad (8)$$

In other words systems with varying gene regulatory network architectures yet identical selection pressures, environment, and phenotype.

Define $V(\tau)$ as the parameterized change of coordinates matrix that preserves B and C , with τ a vector of free parameters. The set of *all* phenotypically invariant (minimal) gene networks is,

$$A(\tau) = V(\tau)A(0)V^{-1}(\tau), \quad (9)$$

and a *Linear Evolutionary System* is,

$$\mathcal{S}(\tau) := \begin{cases} \dot{\kappa}(t) &= A(\tau)\kappa(t) + Bu(t) \\ \phi(t) &= C\kappa(t) \end{cases} \quad (10)$$

That is, all (linear and minimal) mechanisms capable of producing the same phenotype can be realized by a unique choice of τ in $\mathcal{S}(\tau)$.

More generally, we denote by $\mathcal{A}_n(A_0)$ the set of all n -dimensional systems equivalent to A_0 :

$$\begin{aligned} \mathcal{A}_n(A_0) &= \{A : Ce^{At}B = Ce^{A_0t}B \text{ for } t \geq 0\} \\ &= \{A : CA^k B = CA_0^k B \text{ for } 1 \leq k \leq n-1\}. \end{aligned} \quad (11)$$

Equivalence of the two characterizations follows from the Cayley-Hamilton theorem. Usually, the dimension n and the reference system A_0 is implicit and we write only \mathcal{A} .

Regardless of minimality, two systems, even in different dimensions, can have identical external dynamics if they are in \mathcal{A} . This set can be completely parameterized using the *Kalman Decomposition*; the set of all linear gene regulatory networks with equivalent phenotypes can be precisely defined. Note that this implies that there is always more than one possible gene regulatory network architecture per phenotype, however we can analytically describe the complete set.

Example 2 (All Phenotypically Equivalent Cell Cycle Control Networks) *Let*

$$A(0) = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 1 \end{bmatrix}^T, \quad C = \begin{bmatrix} 1 & 0 \end{bmatrix},$$

with $V(\tau)$ preserving both B and C , then the set of all two-gene regulatory networks phenotypically equivalent to the cell cycle control network in Example 1 are given by

$$A(\tau) = \begin{bmatrix} 1 & 0 \\ \tau & 1-\tau \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \frac{\tau}{\tau-1} & \frac{-1}{\tau-1} \end{bmatrix} \quad \forall \tau \neq 1$$

this sentence lacks an object (or some sentence part?)

I don't think it makes sense to write things in general as a function of τ , since in general the set of "free parameters" is complicated and we don't actually specify it. Better, maybe define \mathcal{V} to be the set of all V that preserve B and C , then define $A(V)$ instead of $A(\tau)$...

why is this "evolutionary"? It is the thing that we're thinking about evolution of, but that doesn't make it evolutionary.

again, lacking an object

introduce Kalman

To make sense of this we need to say how C and B change with n .

should say more about this decomposition here.

this is same as above, so reference instead of writing out again?

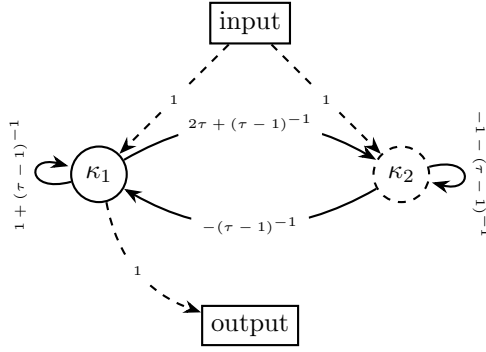


Figure 2: A graphical depiction of the set of all externally equivalent cell cycle control networks, $A(\tau)$. τ can take any value *sentence fragment*

Despite the phenotypic equivalence of all instantiations of $A(\tau)$, the internal dynamics, or kryptotypes, vary as a function of τ . Gene-1 dynamics (blue) are equivalent for network architectures $A(0)$ and $A(2)$, however the dynamics of gene-2 (orange) differ with τ .

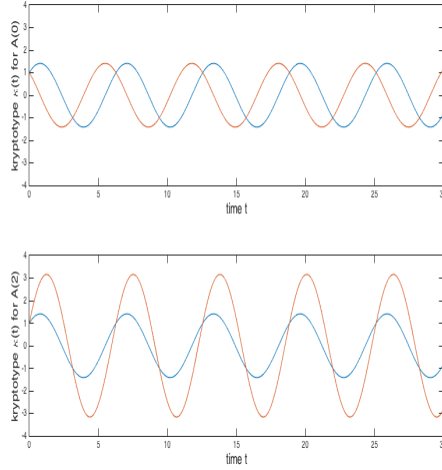


Figure 3: Gene-1 (blue) and gene-2 (orange) dynamics for $A(0)$ (top) and $A(2)$ (bottom). Both (top and bottom) gene-1 dynamics are given by $\kappa_1 = \sin(t) + \cos(t)$, and gene-2 by $\kappa_2 = \cos(t) - \sin(t)$ (top) and $\kappa_2' = \cos(t) + 3\sin(t)$ (bottom).

Missing: a statement about when there are neutral directions, and how many.

I think we should first put the section about quantitative drift in, and then the following section that tries to get at real values of parameters.

Systems drift and the accumulation of incompatibilities

At any given time, there will be a range of regulatory coefficients present in the population due to segregating genetic variants. Over many generations, even if selective pressures do not change, this range of networks will shift as recombination, mutation, and demographic noise create new alleles and shift allele frequencies. How much variation do we expect to find within a population? Is this range limited by available variation or kept in check by selection? How fast will a population explore the space of equivalent networks? In this section we explore

informally a general model for this situation, in which a population drifts stochastically near a set of equivalent, optimal systems. We work with a population of effective size N_e .

Suppose that a set x of coefficients that determine a system (this is A above), produce a phenotype $\Phi(x)$ (the time course of $\phi(t)$). There is an optimal phenotype Φ_0 , and a set \mathcal{X} of “optimal” coefficients that produce this phenotype. Fitness depends on distance to the optimal phenotype – we will write the “distance” between phenotypes ϕ and ψ as $d(\phi, \psi)$, measured so that the fitness of an organism with coefficients x is $\mathcal{F}(x) = \exp(-d(\Phi(x), \Phi_0)^2)$. We will assume that the map Φ is smooth and that the optimal set \mathcal{X} is locally isomorphic to \mathbb{R}^m .

say that better

Offspring. Individuals are diploid; we assume that each haploid genome determines a set of coefficients, and the individual’s coefficients are the average of her two haploid values (no dominance). This implies that the diploid population variance, σ^2 , is one-half the haploid variance. Each new gamete is produced from the parent’s two haploid copies; for simplicity we assume that the gamete inherits a random choice of one of the two parental copies, and so σ remains constant, up to a $1/N_e$ term. A more general model including segregation variance [?] would result in the same qualitative conclusions. *but put this in the appendix?*

System drift If the variation within a population of some coefficient has standard deviation σ , then since subsequent generations resample from this diversity, the population mean coefficient will move a random distance of size $\sigma/\sqrt{N_e}$ per generation, simply because this is the standard deviation of the mean of a random sample [?]. Selection will tend to restrain this motion, but mean movement along the optimal set \mathcal{X} is unconstrained. The amount of variance in particular directions in coefficient space depend on constraints imposed by selection and the covariance between genetic variation between different coefficients (the G matrix [?]). For instance, if the variation is due to *cis*-regulatory variants, the genetic basis of each *row* of A likely lies within a few kilobases of tightly linked sequence, across which a population may carry only a few common haplotypes. However, covariance due to transiently assembled haplotypes is not expected to be stable over long periods of time – a common *cis*-regulatory haplotype of transcription factor k with particularly strong binding to both i and j (leading to positive covariance between A_{ik} and A_{jk}) is no more likely to appear than one with strong binding to i but particularly weak binding to j (negative covariance). (Such transient covariances may well increase the variance of the per-generation change in network mean, however [?].)

To obtain a general quantitative picture, we need to know σ_N and σ_S , the standard deviations of coefficient variation along and perpendicular to \mathcal{X} respectively, and γ , the scale on which phenotype changes moving away from \mathcal{X} . Concretely, γ is the inverse of the derivative of $d(\Phi(x + uz), \Phi(x))$ with respect to u for $x \in \mathcal{X}$ and z perpendicular to the tangent space at x . With these parameters, a typical individual will have a fitness of around $\exp(-(\sigma_S/\gamma)^2)$.

Hybridization By the arguments above, the means of two allopatric populations separated for T generations will be a distance of order $2\sigma_N\sqrt{T/N_e}$ apart along \mathcal{X} . A population of F_1 hybrids has one haploid genome from each, whose coefficients are averaged, and so will have mean system coefficients at the midpoint between their means, and variance equal to σ . Each F_2 hybrid will be homozygous for one parental allele on average at half of the loci in the genome, so the distribution of F_2 s will have mean at the average of the two populations, as before, but variance equal to $\sigma^2 + z^2/2$, where z is the distance between the parental populations. These are depicted in figure 4.

I put F1s and F2s here, and \mathcal{I} , but with less explanation than below. Merge somehow.

improve figure by putting labels on from the following Suppose that two populations have drifted independently to differ by z , and that z is of the same order as σ but is smaller than

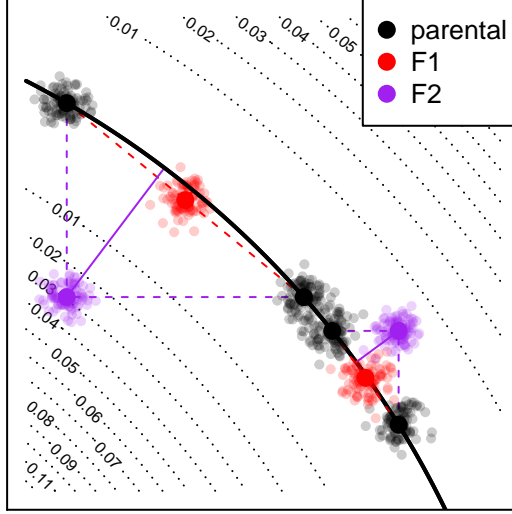


Figure 4: A conceptual figure of the fitness consequences of hybridization: axes represent system coefficients (i.e., entries of A); the line of optimal system coefficients is down in black; dotted lines give phenotypic distances to the optimum. Two pairs of parental populations are shown in black, along the optimum; a hypothetical population of F_1 s are shown for each in red, and the distribution of one type of F_2 is shown in purple (other types of F_2 are not shown). Solid lines depict the distance of the F_2 to optimum. *Should show all types of F_2 ? would be messy.*

γ . The mean F_1 is the average of the parental means, and since the first-order terms in the Taylor series vanish, has phenotype differing from the optimum by a distance of order $\|z\|^2$ (see appendix C). The mean F_2 is the same, but the standard deviation is of order z , so that up to lower order terms, while the typical fitness of an individual in the original population is $\mathcal{F}_0 = \exp(-(\sigma_s/\gamma)^2)$; of an F_1 is $\mathcal{F}_1/\mathcal{F}_0 = \exp(-(c_1\sigma_N^2 T/N_e)^2)$; and of an F_2 is $\mathcal{F}_2/\mathcal{F}_0 = \exp(-T/(N_e\gamma^2))$.

Parameter estimates for real systems

To translate these results into real predictions, we need to know the strength of stabilizing selection on the phenotype, and the amount (and structure) of heritable variation in the genotype. These are known at best only roughly [?], so we aim for order-of-magnitude estimates.

We quantify (roughly) the amount of heritable variation by σ^2 , the genetic variance present in a population in a typical entry of A . The coefficient A_{ij} measures how much the rate of net production of i changes per change in concentration of j . It is generally thought that regulatory sequence change contributes much more to inter- and intraspecific variation than does coding sequence change affecting molecular structure [?]. In the context of transcription factor networks this may be affected not only by the binding strength of molecule j to the promoter region of gene i but also the effects of other transcription factors (e.g., cooperativity) and local chromatin accessibility [?]. For this reason, the mutational target size for variation in A_{ij} may be much larger than the dozens of base pairs typically implicated in the handful of binding sites for transcription factor j of a typical promoter region, and single variants may affect many entries of A simultaneously. On the other hand,

a diverse set of buffering mechanisms are thought to contribute to phenotypic stability in the presence of substantial molecular noise [??], suggesting that substantial variation in the micro-scale dynamics we consider here may be necessary to produce relevant phenotypic effects downstream.

replace "micro-scale" with sthg else or discuss earlier

The amount and structure of this standing variation is established over long time scales by many factors, including mutation-selection balance, shifts in the phenotypic optimum, and/or spatial variation in the optimum [?]. Quantitative genetics models of mutation-selection balance predict precise levels and structure of standing variation [???], but it is unclear how well these predictions match reality [?] and how much they are expected to change over time [?]. However, empirical work allows us to estimate at least the rough magnitude of variation. Differences in A_{ij} due to a sequence change are hard to measure, but variation in both transcription factor binding site occupancy and expression levels (e.g., cis-eQTL) have been measured in various systems. However, variation in binding site occupancy may overestimate variation in A , since it does not capture buffering effects (if for instance only one site of many needs be occupied for transcription to begin), and variation in expression levels measures changes in steady-state concentration (our κ_i) rather than the *rate* of change. Nonetheless, ? found differential occupancy in 7.5% of binding sites of a transcription factor (p65) between human individuals. ? showed that cis-regulatory variation accounts for around 2–6% of expression variation in human blood-derived primary cells, while [?] found that human population variation explained about 3% of expression variation. Taken together, this suggests that variation in the entries of A may be on the order of 1% between individuals of a population – doubtless varying substantially between species and between genes.

Get some data from at least one other species in here!

It seems certain that selection in most species is not so strong that intra-population variation is strongly deleterious, so that if u is the typical scale on which selection acts, then $u > \sigma$. However, a range of studies have found evidence for weak stabilizing selection on regulatory SNPs and cis-eQTL. For instance, ? found evidence that large-effect regulatory mutations are weakly selected against in *Drosophila*. This suggests that the strength of selection on phenotype is sufficient to weakly constrain regulatory variation, so that perhaps σ and u are relatively close. This is as would be expected if available variation is held in check by mutation–selection balance rather than genetic drift. A conservative estimate would be that $u = 5\sigma$; taking $\sigma = .01$ as above, this suggests that changes in phenotype of 5% are sufficient to effect a noticeable drop in fitness.

do we call this β or u ?

find them

(others?)

BUT σ IS VARIATION IN A NOT PHENOTYPE

We have guessed that within a population, the entries of A vary by around 1%, at least for networks whose function is strongly constrained.

Speciation via Reproductive Incompatibility

An F_1 's genetic regulation is a consequence of both of its genomes, where genes and regulatory sequences from both parents are equally present. Thus we say that a diploid organism's gene network is simply the average of both of its gene network copies; one from each parent. Further, each haploid parental gene network copy is formed via meiosis – by swapping independent genes (specifically their *cis*-regulatory modules) randomly. Assuming two distinct but genetically homogeneous populations evolving in allopatry meet and form hybrids, the first generation hybrids (F_1 s) gene network dynamics will be determined by the average of the parental haplotypes. The second generation hybrids (F_2 s), however, will be the product of a meiosis between parental haplotypes followed by an averaging of gametes.

what does "equally present" mean?

should have a stronger justification for the assertion that F_1 is the arithmetic average of the parents than "Thus we say"

Specifically, F_1 is the first generation hybrid gene network architecture formed by mating (averaging) $A(\tau)$ and $A(\hat{\tau})$,

To agree with the previous section we need to incorporate within-pop variation in this argument.

$$F_1(\tau, \hat{\tau}) = \frac{A(\tau) + A(\hat{\tau})}{2}. \quad (12)$$

I don't think we say that A is a "gene network architecture".

and F_2 is the second generation hybrid gene network architecture formed by gametes $G(i)$

mating is not averaging

and $G(j)$. Where each G is formed by randomly swapping rows between $A(\tau)$ and $A(\hat{\tau})$, such that the i th gene comes from $A(\tau)$ (i and j are orthogonal vectors, each element 0 or 1, and $i \neq j$).

sentences don't begin with "Where", even in math. Also, there isn't "the" F_2 - there are many possible.

$$F_2(i, j) = \frac{G(i) + G(j)}{2} \quad (13)$$

The fitness of an organism can be computed by comparing its impulse response with the optimal response,

this is either zero or infinite. apply a weighting function (see appendix)?

$$\mathcal{F}(\hat{\phi}(t)) = \exp \left\{ -\frac{1}{\sigma} \int_0^\infty \left\| \phi(t) - \hat{\phi}(t) \right\|^2 dt \right\}. \quad (14)$$

Therefore a hybrid's fitness can be computed by comparing its impulse response with that of its parents.

Definition 1 (Reproductive Incompatibility) *According to Mayr's Biological Species Concept, two populations can be considered different species if they are reproductively isolated, meaning crosses between them produce low fitness offspring.*

This part doesn't need a definition - we don't need to discuss what is or isn't a "species".

We quantitatively describe the degree of incompatibility between two populations P_1 and P_2 as

$$\mathcal{I} = \frac{2 \langle \mathcal{F}(\phi_{F_1}) \rangle}{\langle \mathcal{F}(\phi_{P_1}) \rangle + \langle \mathcal{F}(\phi_{P_2}) \rangle},$$

where angled brackets imply averaging.

averaging over what? And, this is as others define - cite.

F1s created by crossing phenotypically equivalent oscillators $A(0)$ and $A(2)$ have a phenotype of $\phi_{F_1}(t) = e^t$, in contrast to both parents, who have $\phi_{P_1}(t) = \phi_{P_2}(t) = \sin(t) + \cos(t)$. The hybrid phenotype is significantly different (it does not oscillate and increases infinitely) despite the phenotypic equivalence of the parents.

This is not part of the definition of \mathcal{I} :

$$\begin{aligned} \mathcal{F}(\phi_{P_1}) &= \mathcal{F}(\phi_{P_2}) = 1 \\ \mathcal{F}(\phi_{F_1}) &= 0 \\ \mathcal{I} &= 0 \end{aligned}$$

don't need these to be big equations.

Thus if populations 1 and 2 are homogeneous $A(0)$ and $A(2)$, respectively, we say that they are completely incompatible as $\mathcal{I} = 0$.

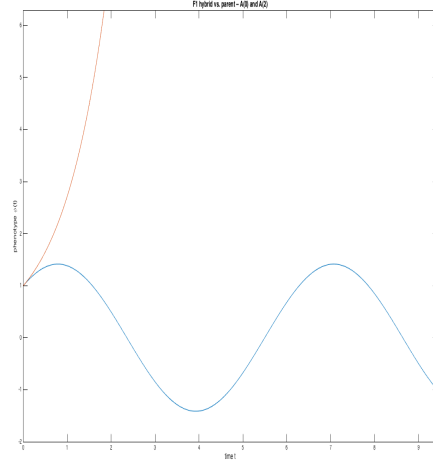


Figure 5: F1 hybrid (orange) and parental (blue) phenotypic oscillator dynamics for an $A(0)$ by $A(2)$ cross. The hybrid fails to oscillate and exhibits qualitatively different dynamics.

Example 3 (Hybrid Incompatibility in an Oscillating Gene Network) *Here we compare the phenotypes for F2 hybrids formed by crossing oscillators $A(2)$ with $A(2.01)$, $A(2.1)$, and $A(2.5)$ (B and C are the same as above). Each A is phenotypically identical ($\phi(t) = \sin(t) + \cos(t)$), however some of the hybrids exhibit markedly different dynamics. These differences tend to increase with time, quadratically (x^2) in F2s, and quartically (x^4) in F1s.*

I don't think this is the right place to introduce this?

$$\begin{aligned}
 A(2) &= \begin{bmatrix} 2 & -1 \\ 5 & -2 \end{bmatrix} & A(2.01) &= \begin{bmatrix} 2 - \frac{1}{101} & -1 + \frac{1}{101} \\ 5 + \frac{1}{99} & -2 + \frac{1}{101} \end{bmatrix} \\
 A(2.1) &= \begin{bmatrix} 2 - \frac{1}{11} & -1 + \frac{1}{11} \\ 5 + \frac{6}{55} & -2 + \frac{1}{11} \end{bmatrix} & A(2.5) &= \begin{bmatrix} 2 - \frac{1}{3} & -1 + \frac{1}{3} \\ 5 + \frac{2}{3} & -2 + \frac{1}{3} \end{bmatrix}
 \end{aligned}$$

F1 gene regulatory networks are formed by averaging the two parental $A(\tau)$ matrices; F2s are formed by first recombining parental matrices followed by an averaging.

cut this sentence, was already said:

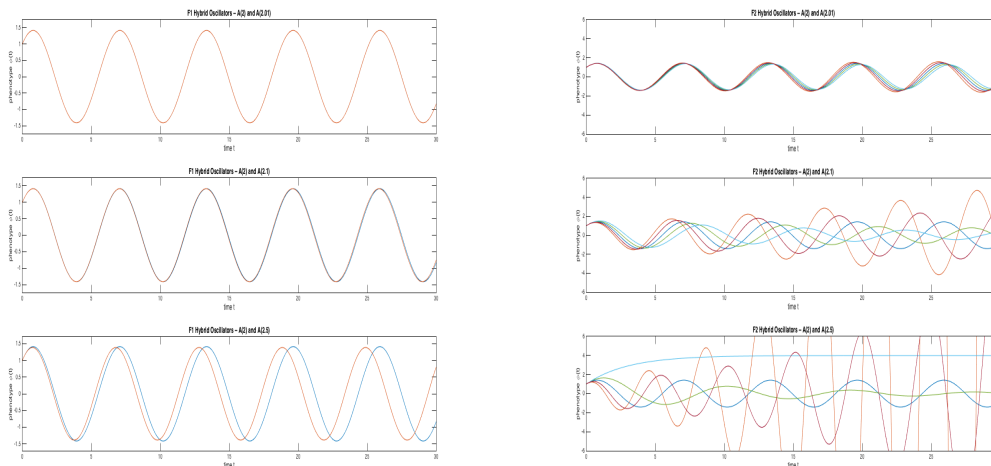


Figure 6: F_1 (left) and F_2 (right) hybrids crossing $A(2)$ with $A(2.01)$ (top), $A(2.1)$ (middle), and $A(2.5)$ (bottom). Note the difference in scale on the y -axis. F_2 hybrids display more phenotypic divergence than F_1 s, on average. Further, some F_2 s completely fail to oscillate, as seen in an $A(2.5)$ F_2 (light blue). *make axis labels bigger*

$A(2)$ and $A(2.1)$ differ in regulatory interaction strengths by 5.09%. If intra-population regulatory variation is approximately 5%, then this level of divergence is expected after $\sqrt{\frac{t}{N_e}} = 1$ generations.

The rate of speciation under neutrality in allopatric populations

In general terms, if fitness \mathcal{F} is some function g of phenotypic divergence from an optimal scaled by the strength of selection β ,

$$\mathcal{F} = g\left(\frac{\phi - \phi_0}{\beta}\right)$$

And a population P is normally distributed around a mean gene network architecture $A(\tau)$ and σ^2 is the intra-population regulatory variance,

$$P_T \sim \mathcal{N}\left(A(\tau), T \frac{\sigma^2}{N_e}\right)$$

the population mean architecture will move at rate

$$\Delta\tau = \sigma \sqrt{\frac{T}{N_e}}$$

where T is time measured in generations. USING GENERAL MATH FROM MULTIVARIATE TAYLOR SERIES... we know that the phenotype of a first generation F_1 hybrid diverges from that of its parents (assuming both parents are at an optimum and phenotypically equivalent) quartically as a function of $\Delta\tau$. Thus the squared difference of it's

separate out here (or above) the difference in genotype of hybrids from the optimum set, the difference in phenotype, and the difference in fitness.

sentences don't begin with "And"

what is T ?

what is $\Delta\tau$? Also, if we reframe in terms of V not τ above will need to reframe this.

phenotype will be,

$$\Delta\Phi_1^2 = c_1^2 \sigma^4 \left(\frac{T}{N_e} \right)^2$$

and F_1 fitness will decline linearly with respect to $\frac{T}{N_e}$. ARGUMENT FOR WHY $\sigma^2 c_1 \sim \beta$ HERE.

$$\Delta\mathcal{F}_1 = \left(\frac{\Delta\Phi_1}{\beta} = \frac{c_1 \sigma^2 \frac{T}{N_e}}{\beta} \approx \frac{T}{N_e} \right)$$

It follows FROM THE SAME MATH ABOUT THE MULTIVAR TAYLOR that the phenotypes of second generation F_2 hybrid crosses diverges quadratically as a function of $\Delta\tau$,

$$\Delta\Phi_2^2 = c_2^2 \left(\sigma \sqrt{\frac{T}{N_e}} \right)^2 = c_2^2 \frac{T}{N_e}$$

$$\Delta\mathcal{F}_2 = \left(\frac{\Delta\Phi_2}{\beta} = \frac{c_2 \sigma \sqrt{\frac{T}{N_e}}}{\beta} \approx \sqrt{\frac{T}{N_e}} \right)$$

and that F_2 fitness drops rapidly – at approximately rate $\sqrt{\frac{T}{N_e}}$. ARGUMENT FOR WHY $c_2 \sigma \sim \beta$ HERE.

This suggests that reproductively isolated populations can speciate rapidly – on timescales much less than N_e generations, depending on the specifics of the fitness function.

Example 4 (Metabolic network) Consider an organism that can metabolize two different sugars s_1 and s_2 (present at logarithmic concentrations u_1 and u_2 respectively), with enzymes e_1 and e_2 (with log concentration denoted as ϕ_1 and ϕ_2). Further suppose that one of the sugars s_2 is the preferred energy source (perhaps it contains significantly more energy than the other sugar or is otherwise more efficiently metabolized). The organism may have a gene regulatory network \mathcal{S} that can deploy a situation specific metabolic strategy. That is depending on both u_2 and u_1 the organism will synthesize an appropriate ϕ_1 and ϕ_2 . Furthermore, consider this system to contain at least two transcription factors, whose log concentrations are given by $\kappa_1, \kappa_2, \dots \kappa_n$.

Minimally such a system may have the architecture,

$$\mathcal{S}_{\min} = \begin{cases} \dot{\kappa}(t) &= \begin{bmatrix} 0 & -1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \kappa_1 \\ \kappa_2 \end{bmatrix} (t) + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} (t) \\ \phi(t) &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \kappa_1 \\ \kappa_2 \end{bmatrix} (t) \end{cases}$$

\mathcal{S}_{\min} is minimal as its reachability and observability matrices both have $\text{rank} = 2$, $\text{rank}(\mathcal{R}) = \text{rank}(\mathcal{O}) = 2 = n \implies \min$.

The impulse response matrix of this system is,

$$h(t) = \begin{bmatrix} 1 & -t \\ 0 & 1 \end{bmatrix}.$$

The phenotype is,

$$\phi(t) = \exp \left\{ \kappa(0) + \int_0^t h(t) u(t) dt \right\}$$

Concentrations are "logarithmic", you mean that $u_1(t) = \log(\text{concentration of molecule 1 at time } t)$?

I think you argue below that the optimal response is a complete switch from one enzyme to the other in the presence of the second sugar? Should say so here.

this is wrong, fix it

($\kappa(0)$ can be set to something like $[-10, -10]$, assuming the transient transcription factor and enzyme concentrations in the organism are typically quite low).

We can see that changing coordinates on this system (with any invertible $V \in \mathbb{R}^{2 \times 2}$) will find new system architectures, as shown before, however, we can also apply the Kalman decomposition to find systems in different dimensions – that is systems that employ more than 2 transcription factors, yet have the same external dynamics/phenotype. This may happen if a system is co-opted from another, or may be a consequence of gene duplication and deletion.

\hat{S} is 4-dimensional biological system,

$$\hat{S} = \left\{ \hat{A} = V \begin{bmatrix} X_1 & X_2 \\ 0 & A_{\min} \end{bmatrix} V^{-1}; \quad \hat{B} = V \begin{bmatrix} X_3 \\ B_{\min} \end{bmatrix}; \quad \hat{C} = [0 \quad C_{\min}] V^{-1} \right\}$$

where V is any invertible 4×4 matrix and X_j is any 2×2 matrix.

So for example, some system can be wired as follows, and still be input-output equivalent to the minimal metabolic system S_{\min} :

$$A' = \begin{bmatrix} 1.6923 & 1.5385 & -2.6154 & -2 \\ 0.8462 & 0.7692 & -1.3077 & -1 \\ 1.2692 & 1.1538 & -1.9615 & -1.5 \\ 0.4231 & 0.3846 & -0.6538 & -0.5 \end{bmatrix}$$

$$B' = \begin{bmatrix} 4 & 4 \\ 2 & 2 \\ 3 & 3 \\ 1 & 3 \end{bmatrix}$$

$$C' = \begin{bmatrix} 0.2692 & 0.1538 & 0.0385 & -0.5 \\ -0.4231 & -0.3846 & 0.6538 & 0.5 \end{bmatrix}$$

Despite the present example consisting of a minimal 2×2 and a non-minimal 4×4 system, any n -dimensional system can be constructed using this method – applying a change of coordinates to the Kalman decomposition – to construct a mechanistically different system with identical phenotypic dynamics. Depending on the specifics of a system being modeled, one may have to take care to restrict the free parameter values and network architectures to be biologically appropriate.

Discussion

Discussion guidelines:

- *Why is this important/useful?*
- *What are the assumptions and shortcomings of the research?*
- *Compare to other studies in the literature.*
- *Future directions.*
- *Wild speculations?*
- *Conclusion and overall impact.*

The complexity of biological systems has limited our understanding of their function and evolution. Above we outline an approach, a first step, towards untangling this complexity in reference to function and evolution. This methodology borrows successfully applied tools

Ah – this is much simpler if we let B and C change as well!

where did \hat{S} come from? need a lead-in here.

the point here is not clear.

what are you thinking of as being a potential problem here? say so explicitly.

mention B part of genetic architecture

from engineering and aims to synthesize these with the concepts and tools of molecular and evolutionary biology.

Theoretical models in evolution and population genetics often lack the molecular details of physiology or of the genotype-phenotype map. Here, we offer a tractable and simple model which includes these missing features. Further, we provide, in clear mathematical language, an analytical description of phenomena hitherto only discussed verbally and conceptually (phenogenetic drift [?], developmental systems drift [?], biological degeneracy [?], *etc.*). The tractability and relative simplicity of this exposition enables the interested biologist to work out by hand, if desired, the dynamics of a genetic system, as well as perturbations to the system – an attribute not likely to be found in less tractable models and simulations.

We have suggested an interpretation of system identification: to see it as an evolutionarily neutral manifold, and not simply a computational nuisance. We have demonstrated a method to analytically determine the set of all phenotypically invariant gene networks; by a simple change of coordinates in the minimal configuration, or more generally by applying the Kalman decomposition in higher dimensions. Further, we emphasize that evolution proceeds through this high dimensional space as stochastic coordinate transformation, constrained by sexual reproduction and selection. This set is explored over evolutionary time when phenotype is conserved, and can lead to a diverse set of consequences, including the accumulation of Dobzhansky-Muller incompatibilities. We emphasize that these incompatibilities are a consequence of recombining different, yet functionally equivalent, mechanisms.

Furthermore, using a quantitative genetic approach, we estimated that a genetically variable population will drift in neutral system space at a rate determined by its intra-population variation and its effective population size. Because mechanistically distinct yet phenotypically equivalent biological systems can fail to produce viable hybrids, we predict allopatric populations to speciate at a rate on the order of N_e under reasonable population genetic parameter estimates. Additionally we see second-generation hybrid fitness plummet much faster than that of first-generation hybrids. This is a consequence of combining our mechanistic model with a quantitative genetic one: we observe that F_1 phenotypes diverge quartically, and F_2 phenotypes quadratically, with evolutionary time. This result is also consistent with Haldane’s rule; that if only one hybrid sex is inviable or sterile it is likely the heterogametic sex. The consistency comes from gene networks localized to the sex chromosomes functioning as an F_2 hybrid cross within a diploid F_1 heterogamete as there is only one sex chromosome.

We also suggest that gene networks may not always use their components parsimoniously as network size tends to ratchet up in the absence of strong selection against extra parts. Although unexplored presently, this phenomena may lead to insights on evolvability and developmental innovation. Lastly, we show that hybrid gene networks break down as function of genetic distance, and may, in part, explain broad patterns of reproductive isolation among diverse phyla [?].

As Richard Levins opined, models in population biology face a trade-off among precision, realism, and generality [?]. As Levins expects, any tractable and general model, such as the present one under discussion, will have limitations. Most notable is linearity. It is often stated that life is not linear. This is often true, however, many of the ideas developed here should be generalizable to nonlinear cases (multi-linear systems, say). Further, we see this as a necessary first step in the direction of more life-like nonlinear evolutionary systems theory. Depending on an actual biological system’s particularities, its (potential) non-linearity, may buffer or exacerbate effects elucidated in this paper, such as the acquisition of Dobzhansky-Muller incompatibilities.

This theoretical framework can easily be applied to other interesting questions in evolutionary biology not tackled presently: such as the evolution of linkage, the necessity of network complexity (does evolution tend towards Rube Goldberg or parsimonious network organization?), evolvability, structure/function inference, and intra-population context dependency of mutational effects, as well as many others.

this is a claim we are the first to do something. best not to make these claims (and we aren’t the first to make math models of these)

probably not desired

check original Bateson/DM papers to see if this accords with those defns

replace "speciate" with "accumulate genetic incompatibilities"

refer to Turelli here

awkward phrasing

probably shouldn’t claim this is "unexplored"

this paragraph isn’t very clear

literature comparison:

Over the last several years, several different computational approaches have been applied to study reproductive incompatibility and speciation. ? simulated the evolution of a transcription factor and its binding site using a thermodynamic model. Their simulations suggest that the language by which a transcription factor recognizes its binding site can change, and potentially lead to hybrid incompatibility when allopatric populations employ divergent readout languages. This study, despite looking at gene regulation, does not analyze overall gene network architecture – as we do here – it only looks at the expression level of a single gene. Furthermore, they report reproductive isolation primarily following directional selection for a change in expression levels in each allopatric population; the evidence for reproductive isolation following balancing selection is much weaker. Johnson and Porter 2000 did not observe any hybrid fitness declines under stabilizing selection – only under directional selection. Khataari et al, Tulchinsky et al, and Porter et al, all study hybrid incompatibility from a transcription factor/binding site interaction perspective, not from an overall network architecture perspective. Palmer and Feldman only see hybrid incompatibility in constant environments if the parental populations are relatively poorly adapted initially. Otherwise hybrids between two allopatric populations have fairly high fitnesses.

citations:

Acknowledgements

We would like to thank Sergey Nuzhdin, Stevan Arnold, Erik Lundgren, and Hossein Asgharian for valuable discussion.

A Kalman Decomposition

Definition 2 (Phenotypic equivalence of systems) Let $(\kappa(t), \phi(t))$ and $(\bar{\kappa}(t), \bar{\phi}(t))$ be the solutions to (??) with coefficient matrices (A, B, C) and $(\bar{A}, \bar{B}, \bar{C})$ respectively, and both $\kappa(0)$ and $\bar{\kappa}(0)$ are zero. The systems defined by (A, B, C) and $(\bar{A}, \bar{B}, \bar{C})$ are **phenotypically equivalent** if

$$\phi(t) = \bar{\phi}(t) \quad \text{for all } t \geq 0.$$

Equivalently, this occurs if and only if

$$h(t) = \bar{h}(t) \quad \text{for all } t \geq 0,$$

where h and \bar{h} are the impulse responses of the two systems.

One way to find other systems equivalent to a given one is by change of coordinates (“algebraic equivalence”): if V is an invertible matrix, then the systems (A, B, C) and (VAV^{-1}, VB, CV^{-1}) have the same dynamics because their transfer functions are equal:

$$CV^{-1}(zI - VAV^{-1})^{-1}VB = CV^{-1}V(zI - A)^{-1}V^{-1}VB = C(zI - A)^{-1}B.$$

However, the converse is not necessarily true: systems can have identical transfer functions without being changes of coordinates of each other. In fact, systems with identical transfer functions can involve interactions between different numbers of molecular species.

The set of all systems phenotypically equivalent to a given system (A, B, C) is elegantly described using the Kalman decomposition, which also clarifies the system dynamics? tells us a lot about how it works? To motivate this, first note that the input $u(t)$ only directly pushes the system in directions lying in the span of the columns of B . As a result, different combinations of input can move the system in any direction that lies in the *reachable*

or something

subspace, which we denote by \mathcal{R} , and is defined to be the closure of $\text{span}(B)$ under applying A (or equivalently, the span of $B, AB, A^2B, \dots, A^{n-1}B$). Analogously to this, we define the *observable subspace*, \mathcal{O} , to be the closure of $\text{span}(C^T)$ under applying A . (Or: $\bar{\mathcal{O}}$ is the largest A -invariant subspace contained in the null space of C ; and \mathcal{R} is the largest A -invariant subspace contained in the image of B .)

If we define

1. The columns of $P_{r\bar{o}}$ are an orthonormal basis for $\mathcal{R} \cap \bar{\mathcal{O}}$.
2. The columns of P_{ro} are an orthonormal basis of the complement of $\mathcal{R} \cap \bar{\mathcal{O}}$ in \mathcal{R} .
3. The columns of $P_{\bar{r}o}$ are an orthonormal basis of the complement of $\mathcal{R} \cap \bar{\mathcal{O}}$ in $\bar{\mathcal{O}}$.
4. The columns of $P_{\bar{r}\bar{o}}$ are an orthonormal basis of the remainder of \mathbb{R}^n .

If we then define

$$P = [P_{r\bar{o}} \mid P_{ro} \mid P_{\bar{r}o} \mid P_{\bar{r}\bar{o}}],$$

then

$$P^T P = \left[\begin{array}{c|c|c|c} I & 0 & 0 & 0 \\ \hline 0 & I & U & 0 \\ \hline 0 & V & I & 0 \\ \hline 0 & 0 & 0 & I \end{array} \right].$$

Check this. Can we get $U = V = 0$?

The following theorem can be found in SOME REFERENCE.

Theorem 1 (Kalman decomposition) *For any system (A, B, C) with corresponding Kalman basis matrix P , the transformed system (PAP^{-1}, PB, CP^{-1}) has the following form:*

$$\hat{A} = PAP^{-1} = \left[\begin{array}{cccc} A_{r\bar{o}} & A_{r\bar{o},ro} & A_{r\bar{o},\bar{r}\bar{o}} & A_{r\bar{o},\bar{r}o} \\ 0 & A_{ro} & 0 & A_{ro,\bar{r}o} \\ 0 & 0 & A_{\bar{r}\bar{o}} & A_{\bar{r}\bar{o},\bar{r}o} \\ 0 & 0 & 0 & A_{\bar{r}o} \end{array} \right],$$

and

$$\hat{B} = PB = \left[\begin{array}{c} B_{r\bar{o}} \\ B_{ro} \\ 0 \\ 0 \end{array} \right],$$

and

$$\hat{C} = CP^{-1} = [0 \quad C_{ro} \quad C_{\bar{r}\bar{o}} \quad 0].$$

The transfer function of both systems is given by

$$H(z) = C_{ro}(zI - A_{ro})^{-1}B_{ro}.$$

In the latter case, we say that the system is *minimal* – there is no equivalent system with a smaller number of species. Note that this says that any two equivalent minimal systems are changes of basis of each other.

Since any system can be put into this form, and once in this form, its transfer function is determined only by C_{ro} , A_{ro} , and B_{ro} , therefore, the set of all equivalent systems are parameterized by the dimension n , the choice of basis (P), the remaining submatrices in \hat{A} , \hat{B} , and \hat{C} (which are unconstrained), and an invertible transformation of $\text{span}(P_{ro})$, which we call T_{ro} .

Theorem 2 (Parameterization of equivalent systems) *Let (A, B, C) be a minimal system.*

- (a) *Every equivalent system is of the form given in Theorem 1, i.e., can be specified by choosing a dimension, n ; submatrices in \hat{A} , \hat{B} , and \hat{C} except for $A_{ro} = A$, $B_{ro} = B$, and $C_{ro} = C$; and choosing an invertible matrix P .*
- (b) *The parameterization is unique if P is furthermore chosen so that each P_x other than P_{ro} is a projection matrix, and that*

$$0 = P_x^T P_y$$

for all (x, y) except (ro, \bar{ro}) .

conjecture:

In some situations we may be interested in only “network rewiring”, where A changes while B and C do not. For instance, if all non-regulatory functions of each molecule are strongly constrained, then C cannot change. Likewise, if responses of each molecule to the external inputs are not changed by evolution, then B does not change.

Another way of saying it: pick the \mathcal{R} and $\bar{\mathcal{O}}$ subspaces, that must intersect in something of the minimal dimension; then let P be the appropriate basis?

A.1 Neutral directions from the Kalman decomposition

The Kalman decomposition above says that any system (A, B, C) can be decomposed into $(P, \hat{A}, \hat{B}, \hat{C})$ so that

$$A = P^{-1} \hat{A} P$$

$$B = P^{-1} \hat{B}$$

$$C = \hat{C} P,$$

and we know precisely how we can change these to preserve the transfer function:

1. $P \rightarrow P + \epsilon Q$ as long as the result is still invertible,
2. $\hat{A} \rightarrow \hat{A} + \epsilon X$ as long as X is zero in the correct places,
3. $\hat{B} \rightarrow \hat{B} + \epsilon Y$ as long as Y is zero in the correct places,
4. $\hat{C} \rightarrow \hat{C} + \epsilon Z$ as long as Z is zero in the correct places.

By taking $\epsilon \rightarrow 0$, these tell us the local directions we can move a system (A, B, C) in. All statements below are up to first order in ϵ , omitting terms of order ϵ^2 .

First, since $(P + \epsilon Q)^{-1} = P^{-1} + \epsilon P^{-1} Q P^{-1}$, modifying $P \rightarrow P + \epsilon Q$ changes

$$\begin{aligned} A &\rightarrow A + \epsilon P^{-1} \hat{A} Q - \epsilon P^{-1} Q P^{-1} \hat{A} P \\ &= A + \epsilon (A P^{-1} Q - P^{-1} Q A), \\ B &\rightarrow B - \epsilon P^{-1} Q B \\ C &\rightarrow C + \epsilon C P^{-1} Q. \end{aligned}$$

Since P is invertible and Q can be anything (if ϵ is small enough), this allows changes in the direction of an arbitrary W :

$$\begin{aligned} A &\rightarrow A + \epsilon (A W - W A), \\ B &\rightarrow B - \epsilon W B \\ C &\rightarrow C + \epsilon C W. \end{aligned}$$

Then, $\hat{A} \rightarrow \hat{A} + \epsilon X$ does

$$A \rightarrow A + \epsilon P^{-1} X P$$

and $\hat{B} \rightarrow B + \epsilon Y$ does

$$B \rightarrow B + \epsilon P^{-1}Y$$

and $\hat{C} \rightarrow C + \epsilon Z$ does

$$C \rightarrow C + \epsilon ZP.$$

These degrees of freedom look like they depend on P , which is not unique, but for any two choices of P there are corresponding choices of X that give the same actual change in A (and likewise for Y and Z).

Therefore, this gives us an upper bound on the number of degrees of freedom, in terms of the dimensions of the blocks in the Kalman decomposition (n_{ro} etc) and the dimensions of B and C (n_B and n_C respectively): namely, for W , X , Y , and Z respectively:

$$n^2 + (n_{r\bar{o}} + n_{ro}n_{\bar{r}o} + n_{\bar{r}\bar{o}}(n_{\bar{r}\bar{o}} + n_{\bar{r}o}) + n_{\bar{r}o}^2) + n_B n_{r\bar{o}} + n_C n_{\bar{r}\bar{o}}.$$

However, some of these may be redundant. For instance, changing P in the direction of a Q that satisfies both $AP^{-1}Q = P^{-1}QA$ and $CP^{-1}Q = 0$ is equivalent to changing B by $Y = QB$.

B Genetic drift with a multivariate trait

For completeness, we provide a brief argument of how the population mean moves under genetic drift with a quantitative genetics model, as in ? or ?. These ignore details of the underlying genetic basis, but developing a more accurate model is beyond the scope of this paper.

Completing the square First note that

$$(x - y)^T A(x - y) = x^T A(x - 2y) + y^T A y,$$

and so

$$\begin{aligned} (x - y)^T A(x - y) + x^T B x &= x^T (A + B) (x - 2(A + B)^{-1} A y) + y^T A y \\ &= (x - (A + B)^{-1} A y)^T (A + B) (x - (A + B)^{-1} A y) + (\text{terms that don't depend on } x). \end{aligned}$$

Therefore, if $f(x; \Sigma, y)$ is the density of a Gaussian with mean y and covariance matrix Σ then substituting $A = \Sigma^{-1}$ and $B = U^{-1}$ above,

$$\frac{f(x; \Sigma, y) f(x; U, 0)}{\int_x f(z; \Sigma, y) f(z; U, 0) dz} = f(x; (\Sigma^{-1} + U^{-1})^{-1}, (\Sigma^{-1} + U^{-1})^{-1} \Sigma^{-1} y).$$

Now suppose that the population is distributed in genotype space as a Gaussian with covariance matrix Σ and mean y . Selection has the effect of multiplying this density by the fitness function and renormalizing, so that if expected fitness of x is proportional to $f(x; U, z)$ then the above argument shows that the next generation will be sampled from a Gaussian distribution with covariance matrix $(\Sigma^{-1} + U^{-1})^{-1}$ and mean $z + (\Sigma^{-1} + U^{-1})^{-1} \Sigma^{-1} (y - z)$. Taking a sample of size N to construct the next generation will produce something close to this but with a slightly (stochastically) deviating mean. The next generation's mean is drawn from a Gaussian distribution with mean with covariance matrix $(\Sigma^{-1} + U^{-1})^{-1}/N$ and mean $z + (\Sigma^{-1} + U^{-1})^{-1} \Sigma^{-1} (y - z)$.

Roughly, what is this doing? Suppose that the population mean differs from the optimum by ϵ , that $\Sigma = \sigma^2 I$ and $U = I/\beta^2$ (so, stabilizing selection happening on a distance scale of β). Then the population mean gets closer to the optimum on average, moving to $\epsilon/(1 + \sigma^2 \beta^2)$ and adds noise of size $(1/\beta)\sigma/\sqrt{N\sigma^2 + N1/\beta^2}$. At equilibrium, these two movements will be of the same order, so that ϵ is of order $(\sigma/\sqrt{N})\sqrt{1 + \sigma^2 \beta^2}$.

C Away from the optimum

Let two points on \mathcal{X} be x_1 and x_2 , let $\bar{x} = (x_1 + x_2)/2$, and let $z = (x_2 - x_1)/2$. Then with $D\Phi$ and $D^2\Phi$ the first and second derivatives of Φ , respectively, then Taylor expanding about x_1 and x_2 finds that

$$\begin{aligned}\Phi(\bar{x}) &= \Phi(x_1) + D\Phi(x_1) \cdot z + \frac{1}{2} z^T D^2\Phi(x_1) z + O(\|z\|^3) \\ &= \Phi(x_2) - D\Phi(x_2) \cdot z + \frac{1}{2} z^T D^2\Phi(x_2) z + O(\|z\|^3).\end{aligned}$$

Now, since $\Phi(x_1) = \Phi(x_2) = \Phi_0$ and

$$\begin{aligned}D\Phi(x_2) &= D\Phi(x_1) + 2z^T D^2\Phi(x_1) + O(\|z\|^2), \quad \text{and} \\ D^2\Phi(x_2) &= D^2\Phi(x_1) + O(\|z\|), \quad \text{and}\end{aligned}$$

adding together the two equations above and dividing by two gets that

$$\Phi(\bar{x}) = \Phi_0 - \frac{3}{2} z^T D^2\Phi(x_1) z + O(\|z\|^3).$$

D Differentiating the fitness function

Suppose that $\rho(t) \geq 0$ is a weighting function on $[0, \infty)$ so that fitness is a function of $L^2(\rho)$ distance of the impulse response from optimal. With A_0 a representative of the optimal set:

$$\begin{aligned}D(A) &:= \int_0^\infty \rho(t) |h_A(t) - h_{A_0}(t)|^2 dt \\ &:= \int_0^\infty \rho(t) |C e^{At} B - C e^{A_0 t} B|^2 dt \\ &= \int_0^\infty \rho(t) |C (e^{At} - e^{A_0 t}) B|^2 dt \\ &= \int_0^\infty \rho(t) C (e^{At} - e^{A_0 t}) B B^T (e^{At} - e^{A_0 t})^T C^T dt\end{aligned}\tag{15}$$

How does this change with A ? Since

$$\frac{d}{du} e^{(A+uZ)t} \Big|_{u=0} = \int_0^t e^{As} Z e^{A(t-s)} ds,\tag{16}$$

we have that

$$\begin{aligned}\frac{d}{du} D(A + uZ) \Big|_{u=0} &= 2 \int_0^\infty \rho(t) C \left(\int_0^t e^{As} Z e^{A(t-s)} ds \right) B B^T (e^{At} - e^{A_0 t})^T C^T dt \\ &= 2 \int_0^\infty \rho(t) C \left(\int_0^t e^{As} Z e^{A(t-s)} ds \right) B (h_A(t) - h_{A_0}(t))^T dt\end{aligned}\tag{17}$$

and, by differentiating this and supposing that A is on the optimal set, i.e., $h_A(t) = h_{A_0}(t)$, (so wolog $A = A_0$):

$$\begin{aligned}\mathcal{H}(Y, Z) &:= \frac{1}{2} \frac{d}{du} \frac{d}{dv} D(A_0 + uY + vZ) \Big|_{u=v=0} \\ &= \int_0^\infty \rho(t) C \left(\int_0^t e^{A_0 s} Y e^{A_0(t-s)} ds \right) B B^T \left(\int_0^t e^{A_0 s} Z e^{A_0(t-s)} ds \right)^T C^T dt.\end{aligned}\tag{18}$$

Here \mathcal{H} is the quadratic form underlying the Hamiltonian. By defining Δ_{ij} to be the matrix with a 1 in the (i, j) th slot and 0 elsewhere, the coefficients of the quadratic form is

$$H_{ij,k\ell}(A) := \mathcal{H}(\Delta_{ij}, \Delta_{k\ell}). \quad (19)$$

We could use this to compute the gradient of D , or to get the quadratic approximation to D near the optimal set. To do so, it'd be nice to have a way to compute the inner integral above. Suppose that we can diagonalize $A = U\Lambda U^{-1}$. Then

$$\int_0^t e^{As} Z e^{A(t-s)} ds = \int_0^t U e^{\Lambda s} U^{-1} Z U e^{\Lambda(t-s)} U^{-1} ds \quad (20)$$

Now, notice that

$$\int_0^t e^{s\lambda_i} e^{(t-s)\lambda_j} ds = \frac{e^{t\lambda_i} - e^{t\lambda_j}}{\lambda_i - \lambda_j}. \quad (21)$$

Therefore, defining

$$X_{ij}(t, Z) = (U^{-1} Z U)_{ij} \frac{e^{t\lambda_i} - e^{t\lambda_j}}{\lambda_i - \lambda_j} \quad (22)$$

moving the U and U^{-1} outside the integral and integrating we get that

$$\int_0^t e^{As} Z e^{A(t-s)} ds = U X(t, Z) U^{-1}. \quad (23)$$

Following on from above, we see that if $Z = \Delta_{k\ell}$, then

$$X_{ij}^{k\ell}(t) = \frac{e^{t\lambda_i} - e^{t\lambda_j}}{\lambda_i - \lambda_j} (U^{-1})_{\cdot k} U_{\ell \cdot}, \quad (24)$$

where $U_{k\cdot}$ is the k th row of U , and so

$$H_{ij,k\ell}(A) = \int_0^\infty \rho(t) C U X^{ij}(t) U^{-1} B B^T (U^{-1})^T X^{k\ell}(t)^T U^T C^T dt. \quad (25)$$

This implies that

$$D(A_0 + \epsilon Z) \approx \epsilon^2 \sum_{ijk\ell} H^{ij,k\ell} Z_{ij} Z_{k\ell} \quad (26)$$

and so

$$D(A_0 + \epsilon Z) \approx \epsilon^2 \sum_{ijk\ell} H^{ij,k\ell} Z_{ij} Z_{k\ell} \quad (27)$$

By section [B](#), if we set $\Sigma = \sigma^2 I$ and $U = H$, then a population at $A_0 + Z$ experiences a restoring force of strength $(I + \sigma^2 H^{-1})^{-1} Z$ (treating Z as a vector and H as an operator on these). If σ^2 is small compared to H^{-1} then this is approximately $-\sigma^2 H^{-1} Z$. This suggests that the population mean follows an Ornstein-Uhlenbeck process, as described (in different terms) in [?](#).