# Reproductive Barriers as a Byproduct of Gene Network Evolution

**Chia-Hung Yang[1] and Samuel V. Scarpino[1,2,3,4,5]***

**\*For correspondence:**
s.scarpino@northeastern.edu

[1]Network Science Institute, Northeastern University, Boston, United States; [2]Department of Marine and Environmental Sciences, Northeastern University, Boston, United States; [3]Department of Physics, Northeastern University, Boston, United States; [4]Department of Health Science, Northeastern University, Boston, United States; [5]ISI Foundation, Turin, Italy

**Abstract** Molecular analyses of closely related taxa have increasingly revealed the importance of higher-order genetic interactions in explaining the observed pattern of reproductive isolation between populations. Indeed, both empirical and theoretical studies have linked the process of speciation to complex genetic interactions. Gene Regulatory Networks (GRNs) capture the inter-dependencies of gene expression and encode information about an individual's phenotype and development at the molecular level. As a result, GRNs can–in principle–evolve via natural selection and play a role in non-selective, evolutionary forces. Here, we develop a network-based model, termed the pathway framework, that considers GRNs as a functional representation of coding sequences. We then simulated the dynamics of GRNs using a simple model that included natural selection, genetic drift, and sexual reproduction and found that reproductive barriers can develop rapidly between allopatric populations experiencing identical selection pressure. Further, we show that alleles involved in reproductive isolation can predate the allopatric separation of populations and that the number of interacting loci involved in genetic incompatibilities, i.e., the order, is often high simply as a by-product of the networked structure of GRNs. Finally, we discuss how results from the pathway framework are consistent with observed empirical patterns for genes putatively involved in post-zygotic isolation. Taken together, this study adds support for the central role of gene networks in speciation and in evolution more broadly.

## Introduction

Over the past 100 years, the role of reproductive isolation due to genetic incompatibilities has received considerable attention in both the empirical and theoretical literature on speciation (*Rieseberg et al., 1996*; *Coyne and Allen Orr, 1998*; *Marques et al., 2019*; *Satokangas et al., 2020*). Through this work, it is widely accepted that divergent selection on *de novo* mutations in geographically isolated populations can facilitate speciation, as originally theorized by (*Bateson, 1909*; *Dobzhansky, 1936*; *Muller, 1942*). Despite well-established examples from *Drosophila* (*Brideau et al., 2006*), *Xiphophorus* (*Wittbrodt et al., 1989*; *Powell et al., 2020*), *Oryza* (*Yamamoto et al., 2010*), *Arabidopsis* (*Bikard et al., 2009*), and *Mus* (*Davies et al., 2016*), the genetics and evolutionary history of incompatibilities are typically far more complex and/or less well understood than what is suggested by classical models (*Noor and Feder, 2006*; *Lowry et al., 2008*; *Presgraves, 2010*; *Wolf et al., 2010*; *Nosil and Schluter, 2011*; *Seehausen et al., 2014*; *Marques et al., 2019*; *Dagilis and Matute, 2020*).

Post-zygotic, genetic isolation is thought to occur due to epistatic interaction between loci, where alleles arise and fix in allopatry prior to secondary contact, e.g., the Bateson-Dobzhansky-Muller

43 (BDM) model (*Bateson, 1909*; *Dobzhansky, 1936*; *Muller, 1942*). However, many incompatibilities
44 uncovered using high-throughput molecular analyses (*Castillo and Barbash, 2017*; *Kuzmin et al.,*
45 *2018*; *Vaid and Laitinen, 2019*) and quantitative trait locus (QTL) mapping (*Moyle and Nakazato,*
46 *2008*; *Turner et al., 2014*; *Chae et al., 2014*; *Lowry et al., 2015*; *Wang et al., 2015*), do not conform
47 to the processes assumed by the BDM model. In particular, in both natural populations and
48 model organisms, studies have found that reproductive barriers often exist between allopatric
49 populations experiencing similar selection pressures and that many of the alleles underlying genetic
50 incompatibility predate the allopatric separation of populations (*Schluter, 2009*; *Han et al., 2017*;
51 *Guerrero and Hahn, 2017*; *Marques et al., 2019*; *Jamie and Meier, 2020*). Both the lack of divergent
52 selection and the role of standing genetic variation are clear violations of the BDM model. As a
53 result, reconciling theoretical models of how and why genetic incompatibilities arise with emperical
54 data on the molecular genetics of post-zygotic, reproductive isolation is of profound importance
55 (*Marques et al., 2019*; *Satokangas et al., 2020*).

56 Analytical and computational models have proposed theoretical explanations for the observed
57 patterns of complex genetic interaction underlying post-zygotic isolation. A collection of models
58 considers *de-novo* mutations at the population level and the accompanying accumulation of hybrid
59 incompatibilities. For example, *Orr* (*1995*) predicted that the number of incompatibilities should
60 increase faster than linearly with the number of substitutions. The study by *Orr* also suggested
61 higher prevalence of complex genetic interactions than simple pairwise incompatibilities. This so-
62 called "snowballing" effect has been further extended by incorporating protein-protein interaction
63 and RNA folding (*Livingstone et al., 2012*; *Kalirad and Azevedo, 2017*). Similarly, *Barton* (*2001*)
64 demonstrated that stabilizing selection can generate hybrid incompatibility between allopatric
65 populations using a quantitative genetics models.

66 The substitution-based approaches, nevertheless, are often at odds with emerging data on the
67 evolutionary history of alleles involved in reproductive isolation (*Marques et al., 2019*; *Satokangas*
68 *et al., 2020*). In addition, many models make an implicit assumption that two allopatric lineages
69 only differ by fixed alleles, which does not capture the empirical diversity among individuals'
70 gene expression (*Kelly et al., 2017*; *Tyler et al., 2017*; *Gould et al., 2018*; *Mogil et al., 2018*; *Ryu*
71 *et al., 2019*) nor the observed importance of regulatory disruption and standing genetic variation in
72 generating reproductive isolation (*Hopkins and Rausher, 2011*; *Guerrero et al., 2016*; *Rougeux et al.,*
73 *2019*; *Morgan et al., 2020*). More importantly, substitutions originating from *de-novo* mutations
74 fail to explain the recent evidence that alleles underlying reproductive barriers often predate
75 speciation events and can evolve along parallel evolutionary trajectories (*Kaeuffer et al., 2012*;
76 *Sicard et al., 2015*; *Meier et al., 2017*; *Nelson and Cresko, 2018*; *Wang et al., 2019*; *Duranton et al.,*
77 *2019*; *Marques et al., 2019*).

78 Another class of computational approaches focuses on the regulation structure that is potentially
79 responsible for complex genetic interactions and resulting incompatibilities. Specifically, researchers
80 consider the evolution of gene regulatory networks (GRNs), which describe the inter-dependencies
81 between gene expression and encode information about both genotype and phenotype. First,
82 *Johnson and Porter* (*2000*) simulated a single linear regulatory pathway as a sequence of matching
83 functions for binding sites, which resulted in reduced hybrid fitness compared to non-epistatic
84 models. Next, *Palmer and Feldman* (*2009*) explored the developmental process where the expres-
85 sion of gene products was iteratively determined through the regulatory networks. Their model
86 demonstrated that, largely as a consequence of the diverse set of possible development pathways,
87 hybrid incompatibilities due to disrupted GRNS could evolve rapidly. More recently, *Schiffman and*
88 *Ralph* (*2018*) modeled gene networks as linear control systems and demonstrated that reproductive
89 isolation can be a consequence of parallel evolution of GRNs with equivalent mechanism. Lastly,
90 *Blanckaert et al.* (*2020*) showed the importance of higher-order interactions and cryptic epistasis
91 for the evolution of reproductive isolation in the presence of gene flow.

92 The implications from these GRN models are not mere outcomes of layering complexity onto
93 existing approaches. Instead, GRNs are a natural extension from lower-dimensional models due to

Manuscript submitted to eLife

94   their close relationship with coding sequences. Ideally, and hypothetically given "omniscience" over
95   the genomes–including comprehension of every fundamental interaction between molecules–one
96   can reconstruct inter-dependencies among genes and obtain GRNs from a bottom-up approach.
97   Of course, this ambition is far from practical and even sounds like a fantasy. Yet, it shows that
98   GRNs are essentially a direct abstraction of the genome sequence. Furthermore, this abstraction is
99   central to the omnigenic perspective of complex traits (*Boyle et al., 2017*). GRNs therefore bridge
100  the gap between inheritance factors and physiological traits, whose dynamics over generations then
101  becomes a candidate for understanding the genetics of speciation due to genetic incompatibilities.
102      To investigate the role of complex genetic interactions in the speciation process, we develop
103  a network-science model for the evolution of GRNs which specifically focuses on the inherited
104  molecular pathways encoded in them. Our approach, termed the pathway framework, considers
105  GRNs as a functional representation of genotype-to-phenotype maps, where proteins are "nodes"
106  in the network and alleles of loci are "edges." Using this framework, we show how a simple model,
107  which includes sexual reproduction, genetic drift, and natural selection, can drive a rapid increase
108  in reproductive isolation between allopatric populations from standing genetic variation under
109  identical selection pressure. Additionally, we find that genetic incompatibilities can frequently
110  involve many loci, i.e., be of higher order, simply as a by-product of GRN evolution. Finally, we
111  conclude the functional redundancy of GRNs is critical for the rapid emergence of reproductive
112  isolation during population divergence.

## Results

### The Pathway Framework: Networks as a Functional Representation of Genetic Interactions

116  Gene interactions networks are conventionally built such that genes are "nodes" and interactions
117  between genes are "edges" or links, for examples see *Tong et al.* (*2004*); *Schlitt and Brazma* (*2007*);
118  *Langfelder and Horvath* (*2008*). Here, we propose an alternative methodology–termed the *pathway*
119  *framework*–for constructing gene interaction networks. The key idea of the pathway framework is
120  to conceptualize genes, or alleles of genes, as "black boxes" that encapsulate how their expression
121  is regulated. More precisely, the pathway framework transforms alleles of genes into directed
122  edges pointing from nodes that are activator/repressor molecules, e.g., transcription factors, and
123  nodes that represent gene products, e.g., proteins. In *Figure 1* we show how: a.) a gene is activated
124  by a transcription factor and generates a protein product (top-right), b.) two genes interact via
125  a transcription factor created by one gene that activates the other (middle-right), and c.) genes
126  can interact via shared transcription factors (bottom-right). As a result of its flexibility, arbitrarily
127  complex genetic interactions can be encoded as "pathways" through a gene interaction network.
128      Importantly, while our proposed representation is closely related to conventional gene interac-
129  tion networks (and a direct mapping between the two always exists when considering interactions
130  mediated by a single class of molecules, e.g., proteins), the pathway framework is often either a
131  more compact and/or informative representation. For example, anytime a gene is regulated by a
132  protein product from another gene, the conventional framework usually includes redundancy that
133  does not appear in the pathway framework, and the pathway framework will capture information
134  not present in the conventional construction, e.g., see Box 1. Because the computational complexity
135  of network analyses often scales non-linearly with the number of edges, switching to the pathway
136  framework can facilitate a more robust exploration of model space.
137      The pathway framework further highlights how phenotypes are a product of both genetics and
138  the environment (not all nodes in the pathway framework need be gene products). Concentrating on
139  the molecular basis of physiological traits, a phenotype can be thought of as the biochemical status
140  of a universal collection of nodes in the pathway framework, e.g., gene products such as proteins or
141  environmental stimuli. Therefore, under the pathway framework, the development of a phenotype
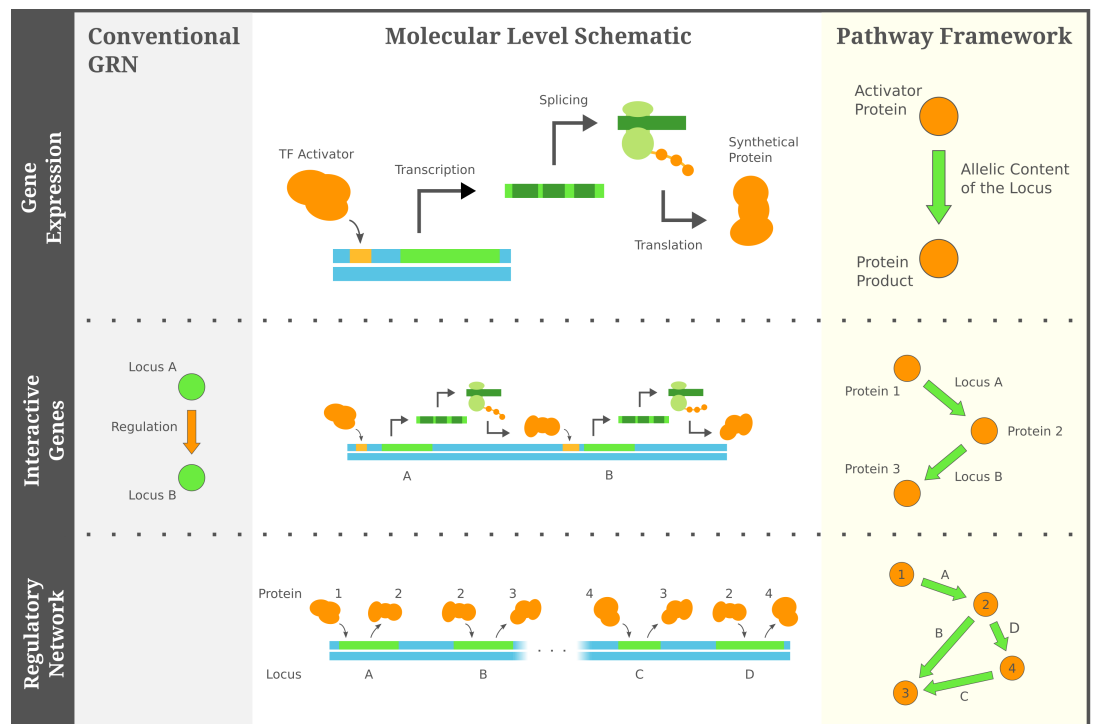142  can be viewed as an iterative process of chemical signals propagating through woven pathways

**Figure 1. The pathway framework captures complex genetic interactions through consecutive regulatory pathways.** In contrast to directly representing genetic interactions, the pathway framework abstracts the regulation and expression of genes as black boxes. If we consider the example of regulation by transcription factors, the pathway framework turns alleles of genes into edges between the transcription factors and the resulting protein products, and regulatory interactions between genes are encapsulated by consecutive pathways through the network.

built from groups of "inherited metabolisms" and external signals from the environment. As a result, the pathway framework can readily capture genetic, environment, and gene x environment effects in the same network.

## Evolutionary Mechanisms under the Pathway Framework

Although in its most abstract state, the pathway framework can include nodes that are not proteins and also nodes that are not directly involved in gene regulation; here, we focus on the evolution of GRNs where all nodes are proteins directly involved in transcriptional regulation. To model and simulate the evolution of GRNs, this version of the pathway framework translates evolutionary mechanisms–such as mutation, independent assortment, recombination, and gene duplication–into graphical operations on the gene networks[1]. Because mutation of a locus can potentially alter its protein product and/or the transcription factor binding region(s), we consider mutation as a rewiring process where the incoming and/or outgoing directed edges are re-directed to point from or to different nodes (*Figure 2*, top-right). Independent assortment during meiosis can be modeled via edge-mixing of parental GRNs such that an offspring acquires alleles, i.e., edges in the GRN, from both parents (*Figure 2*, bottom). Similar to mutation, recombination is an edge-rewiring process that is constrained to swapping binding sites or transcription factors at the same locus. Finally, gene duplication is equivalent to adding a parallel edge that represents the identical allelic content of a duplicated locus.

An individual's viability subjected to natural selection is a response to its molecular phenotypic status, which–under the pathway framework–can be modeled as a fitness function associated with

---

[1]These graphical operations focus on edges in the GRNs, while the underlying node set is held constant because the nodes represent all *possibly existing* proteins in the organism.
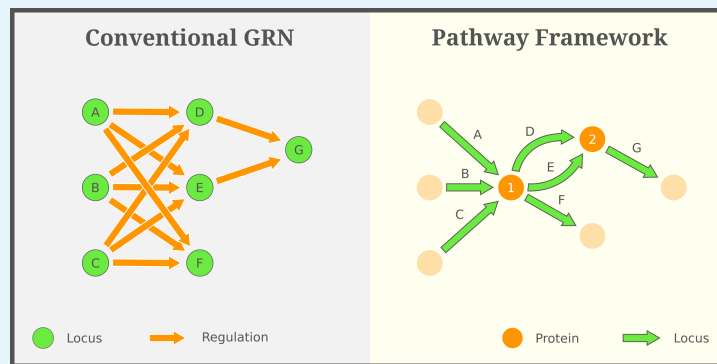
## Box 1. The pathway framework is often a more compact representation

Because the pathway framework directly encodes the expression pattern of genes, it can contain more information than the "conventional" approach to constructing GRNs. When considering genetic interactions that are mediated by a single class of molecules, e.g., one gene being regulated by the protein product of another, the pathway framework takes advantages of this information and presents genetic interactions in a more compact format. Conversely, a conventional GRN lacks the specific regulatory context, and thus it has to present all pairs of interacting genes as individual edges, rather than summarizing these interactions by a smaller set of protein mediators. More technically, the pathway framework and a conventional GRN correspond to the first- and second-order de Bruijn graph (*De Bruijn, 1946*) respectively, where higher orders usually introduce redundant elements and additional computational complexity.



the collective state of nodes and edges in the GRN. For example, one could study the time-varying concentration of each protein, attach a continuous dynamic or a stochastic reaction to every allele and define fitness as a function of the high-dimensional concentration vector, etc.. On the other extreme, we can consider Boolean networks, which have been shown to effectively capture many of the most relevant dynamical features of empirical regulatory systems (*Davidich and Bornholdt, 2008*). In this minimal scenario, each protein is assigned to a Boolean state (present or absent) and external environmental signals stimulate the existence of specific proteins in the organism. The logical states then cascade through the genetic pathways, where–given the presence of a gene's transcription factor–loci activate and generates protein product(s). The phenotype of a GRN is thus the "reachability" from the environmental stimuli, whose binary survival is defined via a sharp fitness landscape over plausible collective Boolean states (*Figure 2*, top-left).

We adopt the Boolean-state assumption of GRNs because they readily shed light on the formation of hybrid incompatibilities. Hybrid incompatibilities are lethal combinations of alleles that were not prevalent or present in parental lineages, but are in hybrids. Moreover, the combination is minimal in the sense that the lack of any of its allelic elements will not lead to an inviable hybrid. In the pathway framework, suppose that binary viability only depends on a set of lethal proteins, i.e. an individual will not survive selection if any of those protein are present, a combination of alleles that includes a pathway from a environmental stimulus to a lethal protein makes the GRN inviable. If the alleles exactly comprise a simple path, which contains no cycles, they become a minimal combination and thus form an incompatibility. Additionally, The complexity of genetic interactions can be characterized by the number of alleles involved, which is called the order of hybrid incompatibility and related to the length of the simple pathway[2].

---

[2]Since for $n \geq 1$, $n + 1$ alleles form an $n$th-order incompatibility, the order of genetic interaction is then the path length minus one.
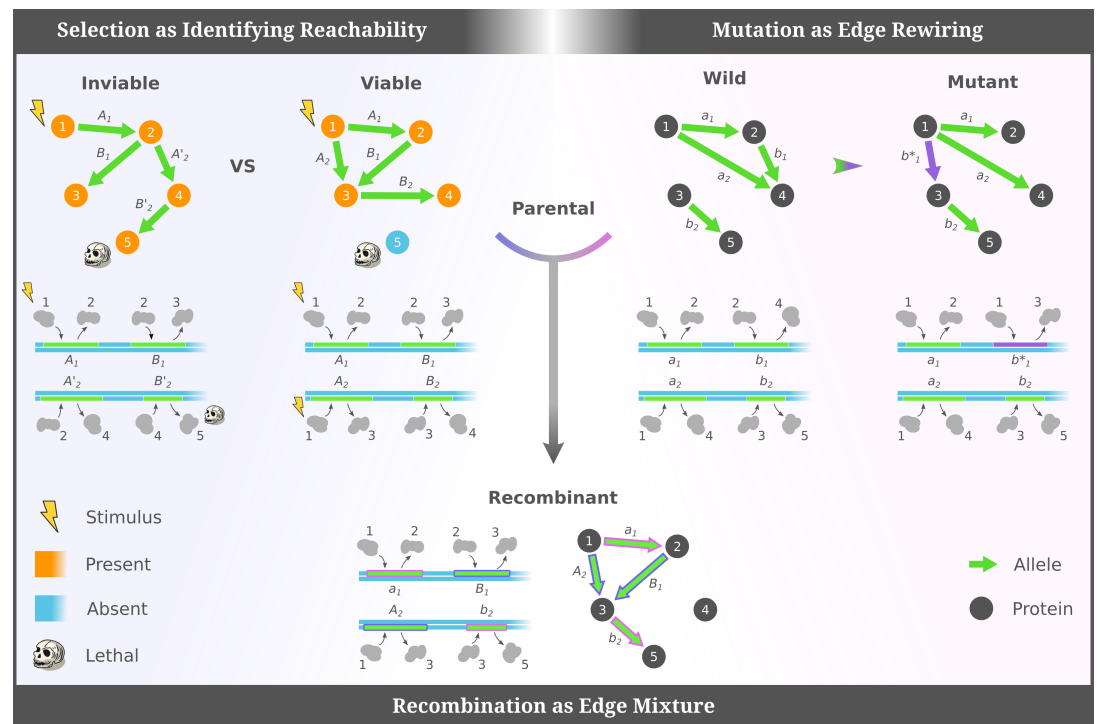
**Figure 2. How the pathway framework turns evolutionary mechanisms into graphical operations on the GRNs.** Since the pathway framework directly models the functionality of alleles of genes as edges, mutation, meiosis, and recombination can be modeled as edge-rewiring and edge-mixing, while a minimal selection scenario of binary fitness can be modeled as identifying "reachability" in a GRN.

## Simulating the Evolution of GRNs

Briefly, we consider a Wright-Fisher model of evolution with natural selection, i.e., constant population size, no mutation, no migration, non-overlapping generations, and random mating. Selection occurs during the haploid stage of the life-cycle, where individuals that survive selection fuse randomly, i.e., create diploids, and undergo meiosis to generate the subsequent generation. Populations are seeded such that each individual has a randomly generated GRN and evolve until a single GRN fixes in the population. Simulations are further detailed in the Methods.

*Figure 3a* shows the proportion of individuals in the population that survive natural selection. Initially, due to the variation of randomly seeded GRNs, the fraction of viable individuals differed substantially between simulations with different initial conditions. However, as the gene networks evolved, the population's viability increased and quickly reached a state where every individual survived selection (dashed line). During this 100% survival stage, natural selection was no longer effective and the population evolved to fixation via genetic drift. Not surprisingly, our results demonstrate that GRNs can rapidly evolve from a heterogeneous population with low average viability to "match" an imposed selective regime or environment.

In addition to achieving 100% survival, populations always fixed a single GRN. *Figure 3b* plots the number of structurally-distinct GRNs in each generation. The decreasing trend demonstrates that, although various GRNs have equal survival probability, it becomes more and more likely that individuals shared a common GRN. Moreover, the populations always fixed a single GRN (dotted line) after a sufficiently long period of time. This phenomenon can be intuitively explained by the mechanism of sexual reproduction. In our model, parents with identical GRNs would lead to offspring of the same GRN, since any two corresponding groups of segregated alleles retrieved the parental gene network. Thus once there was a majority GRN in the population, it would have a higher chance of retaining its genetic configuration in the next generation, as compared to being
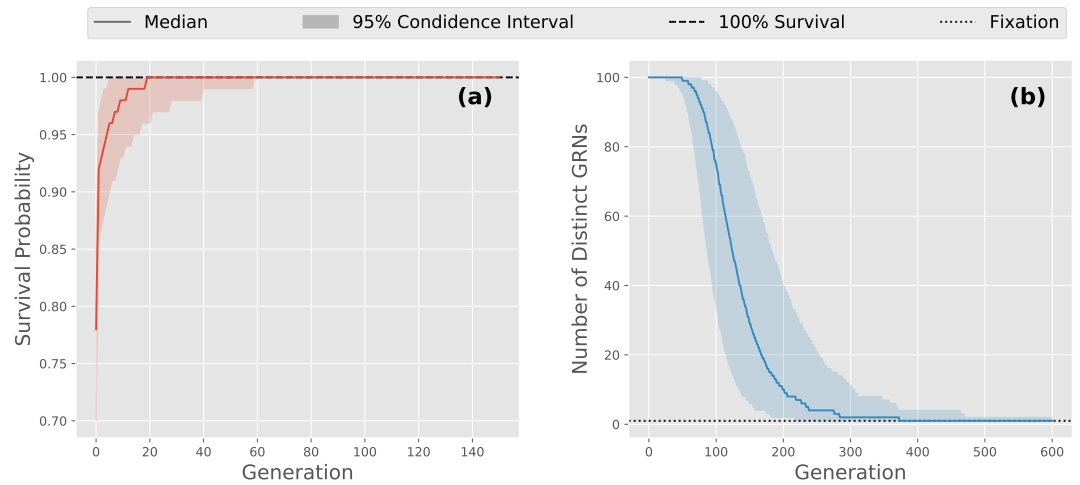
**Figure 3. Populations adapt to the environment and then fix a single GRN.** Here, we show for every generation of GRN evolution, across multiple allopatric populations with different initial conditions: **(a)** the survival probability of an individual and **(b)** the number distinct GRNs in each population, where two individuals' GRNs were deemed effectively identical if they were isomorphic. The average viability of each population increased over time and rapidly achieved 100% survival, which indicates that evolution of GRNs drove adaptation toward the imposed environment. We also observe decreased variation of GRNs as they evolved, with individuals in the same allopatric population, i.e., simulation run, eventually fixing for the same GRN.
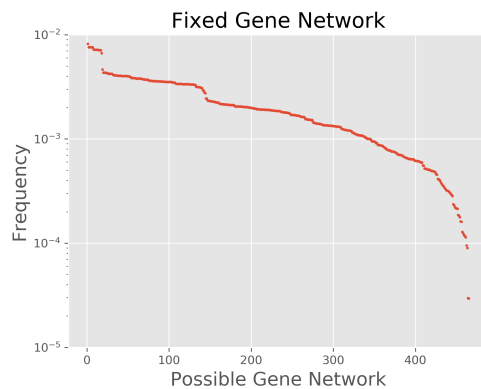


**Figure 4. Fixation of parallel lineages resulted in a wide range of GRN structures.** We simulated isolated populations from the same initial conditions until they reached fixation. In this case Setup 2 in Methods was applied in order to tractably enumerate all plausible GRN, and the ancestral populations were chosen such that the fixation was unbiased by the initial allele frequencies. The $10^7$ acquired GRNs were categorized into 465 viable structures and the fixation frequency of each structure was plotted in a descending order. The distribution shows that isolated lineages fixed alternatives gene networks, some among which were more favorable under our model of GRN evolution.

replaced by meiotically shuffled variants.

Lastly, to better understand how parallel lineages evolve, we consider a scenario where multiple allopatric populations are seeded with the *same* initial conditions. Similarly, each allopatric population rapidly achieved 100% survival and then fixed a single GRN. However, across allopatric populations seeded from the same initial conditions, many different GRNs fixed. *Figure 4* presents the distribution of fixed GRNs for a smaller-scale simulation (Setup 2 in Methods). We see that the fixed GRNs were diverse and non-uniformly distributed. Despite being under identical selection forces and having the same initial condition, lineages evolving from a common ancestral population fixed alternative GRNs. This result demonstrates that a broad range of GRNs can survive the given selection pressure. Furthermore, none of the viable GRN structures had a zero fixation probability, indicating a thorough exploration of evolution in the space of possible GRNs. That so many different GRNs fixed suggests that evolution was less governed by a definite trajectory, but instead it occurs via an uncertain realization among all the possibilities constrained by the ancestral population and the selection pressure.

### Reproductive Barriers Arose Rapidly as Gene Networks Evolved

If the survival probability and fitness of GRNs were identical, the distribution of fixed networks should be uniform over all viable conformations. Because we observe a strongly non-uniform distribution (see *Figure 4*), some other form of selection, i.e., as opposed to simply viability selection, is likely operating on the GRNs. We note that during random mating, even between two parents with viable GRNs, some of their shuffled offspring can be inviable. Coupled with the observation that different allopatric populations, i.e., simulation runs, fix alternative GRNs from the same initial conditions, we hypothesized that some degree of reproductive isolation may exist between these fixed populations.

To test for the presence of reproductive isolation, we performed a "hybridization" experiment between parallel lineages that had reached fixation. Starting with lineages branched from a common ancestral population, two fixed lineages were randomly selected and interbred. Hybrids were generated and the reproductive isolation metric (RI) between the parental populations was computed (see Methods). By repeating this procedure, we obtained a distribution of reproductive isolation, as demonstrated in *Figure 5a* inset. Despite a large fraction of crosses resulting in nearly zero RI, we discovered pairs of lineages with positive reproductive isolation metric. Specifically, the RI distribution displays several regions of positive reproductive isolation such that a high percentage of hybrid offspring are inviable. Thus, we conclude that reproductive barriers between fixed lineages, derived from the same initial population and experiencing identical selection, exist.

Given noticeable reproductive barriers between fixed lineages, we further studied when those barriers first manifested during GRN evolution. Note that because our simulations did not contain mutation, incompatibilities arise because of shuffling during meiosis. Here, instead of waiting until GRN fixation, we instead evolve lineages for $T$ generations and then cross them to generate hybrids as described above. By varying $T$, a series of reproductive isolation distributions were acquired. *Figure 5a* collects and displays them in a heat map. A vertical slice represents a RI distribution as in the inset panel, but crosses were made after $T$ generations rather than waiting for lineages to reach fixation. We see that the regions of high incompatibility noted in *Figure 5a* inset becomes bands in the heat map, which allows us to trace the emergence of reproductive barriers.

Initially the reproductive isolation distribution was relatively symmetric around zero. However, As GRNs evolved, the range of RI broadened and its extreme value in the positive tail increased. The trend towards higher levels of RI decelerated after 100 generations; it then stabilized and formed a band structure, where crosses cluster around certain levels of reproductive isolation. *Figure 5a* hence reflects that reproductive barriers existed at low levels as soon as the lineages started evolving independently and peaked at a time prior to GRN fixation. By assumption, the alleles underlying RI were present in the ancestral population, but we further conclude that RI peaked well before fixation of GRNs.

Next, for incompatible hybrids generated in our crossing experiment, we determine how complex the underlying mechanism of RI was. Specifically, *Figure 5b* shows how frequently an inviable hybrid resulted from an incompatibility of a certain order. We see that hybrid incompatibilities spanned a broad range of interaction orders. Importantly, the simple two-allele interaction was only slightly more common than incompatibilities resulting from three or four interacting alleles and interactions above forth order made up almost 3% percent of all incompatibilities. However, we note that the frequencies of incompatibility order varied depending on the ancestral population.

The pattern of complex genetic interactions provides insights on the distribution of reproductive isolation. Based on the independent assortment mechanism in our model–and assuming that multiple incompatibilities rarely occurred between two parental GRNs–we conclude that hybrid incompatibilities quite often involved higher order interactions, which did not arise as a result of selection, but simply were an expected consequence of GRNs being high order (*Appendix 1*). Further, the discrete characteristic of hybrid incompatibilities led to a higher likelihood at certain RI levels. The band structure in *Figure 5a* agrees with this prediction (*Appendix 1*), which suggests
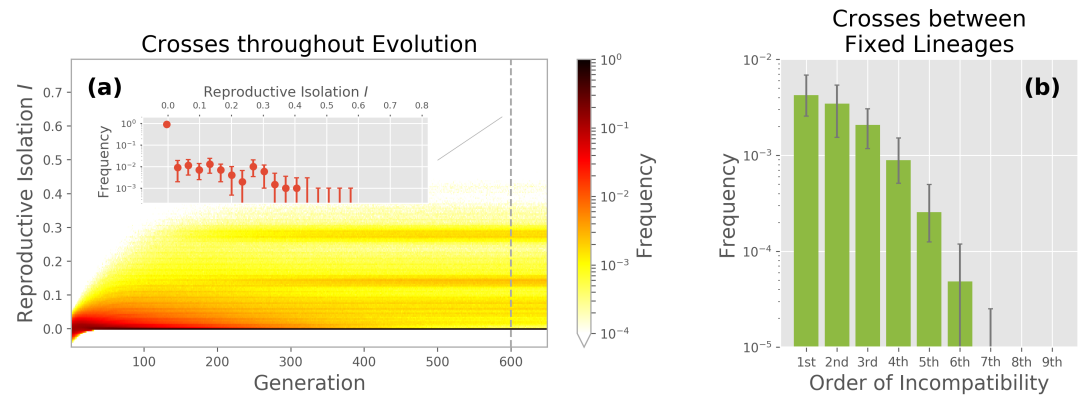
Manuscript submitted to eLife



**Figure 5. Reproductive barriers arose rapidly between allopatric populations. (a, Inset)** Distribution of reproductive isolation between pairs of fixed lineages. A non-negligible fraction of crosses led to positive reproductive isolation, which reflects the occurrence of inviable hybrids and indicates reproductive barriers between fixed lineages. **(a)** We crossed allopatric populations at every generation during GRN evolution and stacked the RI distributions into a heat map. A vertical slice in this heat map represents the RI distribution at a given time, similar to the inset, but where the color shows the mean frequency for each bin. The growing level of positive RI indicates that reproductive barriers arose at the early stage of evolution. **(b)** Frequency that incompatibilities with various order were observed among hybrids between fixed lineages. We see that the order of incompatibilities included a broad range and that the simple pairwise interaction did not significantly dominate over more complex incompatibilities. Moreover, hybrid incompatibilities are consistent with the clustered level of RI and hence sheds light on the observed RI distribution (*Appendix 1*). In both the inset and panel (b), the plots show the statistic of the distribution among multiple groups of allopatric populations, specifically the median frequency and the 95% confidence interval.

287    that reproductive barriers are strongly influenced by the concealed hybrid incompatibilities and are
288    coupled with the genetic interaction pattern shown in *Figure 5b*.

### Early Divergence between Lineages was Critical for Reproductive Barriers to Emerge

290    To further study the emergence of reproductive barriers in our model, we investigated the relative
291    importance of various evolutionary forces in generating the observed patterns of RI. In particular,
292    were the barriers attributed to selection pressure, random genetic drift, or both? We designed
293    two "control scenarios" that were based upon the previously simulated model, but contained
294    modifications to remove the effects of either selection or drift. Comparing the strength and pattern
295    of RI resulting from the two control scenarios, i.e., the removal of drift or selection, to the original
296    GRN dynamics, which contain both evolutionary forces, provides an assessment of the removed
297    component's role in shaping the observed pattern of RI.

298    Removing the effect of natural selection is straightforward to simulate. In this control scenario,
299    populations simply evolve in a selectively neutral environment where all GRNs are viable. Thus, all
300    individuals survived and genetic drift became the only remaining evolutionary force. Of course,
301    this neutrality concurrently made the RI metric ill-defined. We avoided this issue in the crossing
302    experiments to calculate RI by placing the parental populations under the same non-neutral
303    environment in the original model, so the hybrids would be generated from survivors subjected to
304    selection pressure. The reproductive isolation metric could then be computed with respect to the
305    non-neutral environment. Placing the parental population through a round of viability selection
306    just prior to hybridization ensures comparability between the model and the "no selection" control
307    scenario since the survivability of hybrids was evaluated under the same environment and was not
308    biased by the otherwise inviable parents.

309    *Figure 6a* shows the contrast of barriers observed in the original GRN evolution model (red) and
310    in the scenario with no selection (blue). We traced a measure of reproductive isolation over time,
311    defined as the 99th percentile of the RI distribution, which is a sufficient indicator of reproductive
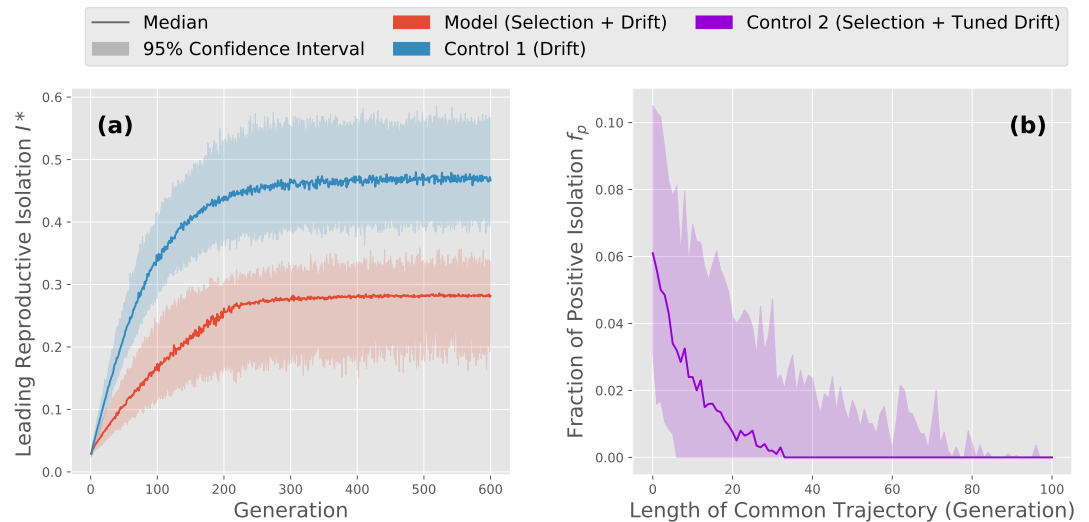
**Figure 6. Early divergence of evolutionary trajectories between lineages was necessary for reproductive barriers to arise.** Here we compare a statistic, termed leading reproductive isolation $I^*$ (99th percentile of the RI distribution), measuring the degree of reproductive barrier in the original model and two designed control scenarios. Control scenarios were simulated with the same group of ancestral populations as the model, where lineages were then crossed to generate hybrids. **(a)** Leading reproductive isolation $I^*$ among allopatric populations over time, where positive values indicate the existence of reproductive barriers. We plot the original model in red and the control scenario with a neutral environment in blue. The increasing and larger $I^*$ uncovered in the control scenario implies that reproductive barriers were still observed when the selection forces were silenced. **(b)** Long-term fraction of positive RI $f_p$ when the influence of random genetic drift was tuned. We simulated the evolution of lineages, but first confine them to a common trajectory of length $L$, which was realized by evolving a single population from the ancestors for $L$ generations, and then simulated allopatric evolution from this now less diverse ancestral population. The original model corresponds to the case where $L = 0$, and for any positive $L$ the effect of drift were lessened. We obtained the $f_p$ metric when lineages evolved for 600 generations, where $f_p = 0$ suggests no barriers among populations. That $f_p$ decreased with $L$ to 0 shows that reducing the effect of drift diminished reproductive barriers. As a result, it implies the criticality of divergence among evolutionary trajectories for barriers to emerge.

barriers between lineages. We discovered that in both the model and the control scenario, the leading RI $I^*$ increased and then saturated. Furthermore, the growth in $I^*$ decelerated after a similar number of generations in both scenarios. That RI occurs at a higher level in the control experiment indicates that selection did not "cause" the fixation of barriers between allopatric populations, but instead suggests that selection was actually limiting chances for incompatibilities to occur in hybrids. We hypothesize that–although restricted as compared to drift–selection operating on incompatibilities likely induced the observed disconnect between viability and fitness seen in *Figure 4*.

We next turned to the contribution of genetic drift. This control scenario, however, was less straightforward due to technical difficulties associated with directly removing random genetic drift from the model. Neither abandoning sexual reproduction nor simulating an infinite population would result in non-trivial and/or computationally tractable GRN evolution. Alternatively, we designed a control scenario where the evolutionary influence of drift could be tuned and limited. Genetic drift results in stochasticity and causes populations to experience diverse trajectories. On the other side of the coin, if two lineages show similar evolutionary trajectories, one would say that drift effectively leads to less divergence between them. We restricted the influence of genetic drift by first confining lineages in a common trajectory for $L$ generations, and then freed the populations and let them evolve independently, i.e., in allopatry. Varying the length of the common trajectory $L$ tunes the overall similarity among lineages. Therefore, $L$ quantitatively reflects the strength of genetic drift.

Manuscript submitted to eLife

*Figure 6b* demonstrates the long-term fraction of positive reproductive isolation introduced in Methods, termed $f_p$, as we varied the length of the common trajectory. Despite substantial variation in $f_p$ in the original model, which corresponds to the case where $L = 0$, a decline of $f_p$ was uncovered as early evolutionary confinement was extended. We discovered 50% of the experiments showed a zero $f_p$ after lineages were evolved together for 40 generations, and as the length of common trajectory exceeded 80 generations positive reproductive isolation was hardly found between lineages. More importantly, *Figure 6b* suggests that as the evolutionary influence of genetic drift was mitigated, RI was weakened and eventually vanished. Namely, restricting early divergence among populations due to genetic drift diminished reproductive barriers. This control scenario consequently suggests that divergence between lineages, coupled with high diversity in the ancestral population, is critical for reproductive barriers to arise.

## Intra-lineage Incompatibilities were Eliminated Stochastically While Inter-lineage Incompatibilities Persisted and Led to Reproductive Barriers

To better understand how reproductive barriers might be removed within a lineage, but persist between lineages, we computed two quantities from the underlying genetic pool. First, the size of the genetic pool, which determines how many possible genotypes a population contains. This measure captures the potential genetic diversity in the population. Second, we count the number potential incompatibilities in the underlying genetic pool, which are lethal allelic combinations that could potentially be realized in the next generation. These incompatibilities compose the source of inviable offspring and RI between allopatric populations. However, because even for small GRNs searching for all possible incompatibilities quickly becomes computationally intractable, we developed a novel algorithm (summarized in Methods) to compute their number in the genetic pool.

Because our model does not contain mutation, one would expect the size of the underlying genetic pool to decline in our simulated gene network evolution. Any allele in an individual was inherited from its parents, and thus it must appear in the parental generation as well. Additionally, a parental allele might not persist in the offspring for two possibilities: either it was not transmitted because of finite population size of the progeny generation and the stochasticity during sexual reproduction, i.e. drift, or it formed a lethal pathway along with other inherited alleles which made the offspring inviable, i.e. selection.

*Figure 7a* demonstrates the size of genetic pool over time, where we compare simulations in the original model (red) and in the control scenario without selection pressure, i.e., only genetic drift will reduce the size of the genetic pool (blue). A rapid decline of genotypic diversity was witnessed under both models. More intriguingly, little difference was found between the GRN evolution model and the control scenario under a neutral environment. The two median curves nearly overlaps, and for any given generation, the pool size in the original model was not significantly smaller than the control counterpart. Therefore, we find additional support for our earlier finding that although both natural selection and random genetic drift decreased genotypic diversity, drift was the dominant driving force. However, while the effect of drift reduced diversity within a lineage, it increased divergence among lineages.

*Figure 7b* shows the number of potential incompatibilities within a lineage's underlying genetic pool (orange). We found that the amount of incompatibilities embedded in a population also decreased over time. This phenomenon is understood by the continual loss of allelic diversity, since removing an allele from the underlying pool always restricts the possibilities to form a lethal pathway in the GRN. Furthermore, the number of potential incompatibilities fell rapidly until no potential incompatibilities remained. The elimination of potential incompatibilities illuminates how a population adapted to the imposed environment when GRNs evolved, as shown in *Figure 3a*. Random genetic drift drove the loss of a lineage's genotypic diversity, and along with the guidance of selection, it eliminated probable lethal pathways in the genetic background. Once all the potential incompatibilities were eliminated, no source of inviable offspring existed and consequently the
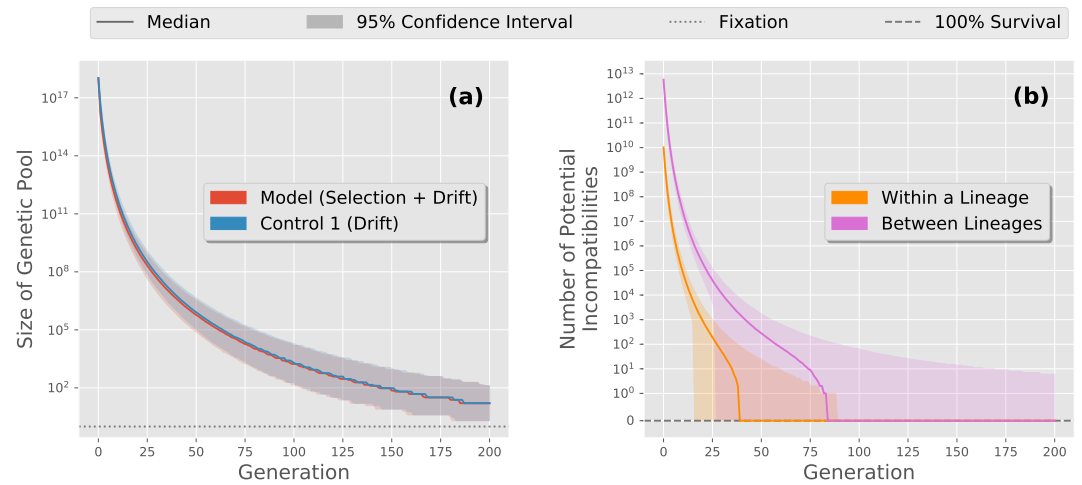
Manuscript submitted to eLife



**Figure 7. The underlying genetic pool lost alleles and eliminated potential incompatibilities within allopatric populations, whereas inter-lineage incompatibilities persisted. (a)** Size of the underlying genetic pool for each generation, where we plot the original model in red along with the no selection control scenario in blue. Both cases show a similar reduction in the genetic pool. The similarity of these curves suggests that the continual losses of allelic diversity within a lineage was dominated by random genetic drift. **(b)** Number of potential intra-lineage (orange) and inter-lineage (pink) incompatibilities for each generation in the original model. We found that the number of potential incompatibilities also decreased as GRNs evolved, which is explained by the reduced allelic diversity in the genetic background. The vanishing intra-lineage incompatibilities implies disappearing sources of inviable hybrids, and it provides a mechanistic understanding of how a genopytically rich populations adapted to the imposed environment. Contrarily, the intra-lineage incompatibilities remained during GRN evolution. It was the persistent potential incompatibilities between allopatric populations that led to evident reproductive barriers.

**Figure 7–Figure supplement 1.** Inter-lineage incompatibilities were sustained throughout GRN evolution.

---

population reached 100% survival. Again, this result supports our earlier finding that natural selection was operating against incompatibilities within a lineage, but that drift was nevertheless the dominate force in structuring incompatibilities between lineages.

Finally, we investigated incompatibilities between underlying pools of lineages, which we call the "inter-lineage" incompatibilities, as compared to potential lethal allelic combinations within a population termed "intra-lineage" incompatibilities. *Figure 7b* presents the number of inter-lineage incompatibilities over generations (pink). We observed more incompatibilities between allopatric populations than those within a population, and similarly their amount dropped as allelic diversity decreased. In contrast, inter-lineage incompatibilities were removed at a slower pace compared to intra-lineage incompatibilities. The sustained confidence interval further suggests that some inter-lineage incompatibilities persisted, which was also the case after populations reached fixation (*Figure 7–Figure Supplement 1*). The persistence of these potential incompatibilities qualitatively explain the inviable hybrids revealed after GRN evolution. In spite of lineages adapting to the same imposed environment, hybrdiziation can "resurrect" a lethal combination of alleles, which was eliminated in either lineages yet remained in their joint genetic background. This explanation also supports the stronger barriers uncovered in the neutrally evolving control in *Figure 6a*, since inter-lineage incompatibilities would be more persistent without the constant selection pressure (*Figure 7–Figure Supplement 1*).

## Discussion

In this work, we develop a pathway-oriented construction of GRNs where alleles are represented as edges in a network. Termed the pathway framework, this model allows us to apply network science analyses to the study of speciation. Specifically, we simulate the evolutionary dynamics of GRNs under a model that includes natural selection, sexual reproduction, and genetic drift. Starting

Manuscript submitted to eLife

from a diverse ancestral population, we show how reproductive isolation can arise rapidly between allopatric populations experiencing identical selection pressure. Then, using a series of counter-factual simulations, we disentangle the relative importance of each evolutionary force included in our model and identify the central roles of high-dimensionality and functional redundancy, even in comparatively small GRNs, for speciation. Finally, we show how higher-order genetic incompatibilities can often evolve simply as a by-product of GRN evolution.

Our counter-factual simulations reveal that the observed reproductive barriers likely resulted from divergent evolutionary trajectories and persistent, inter-lineage incompatibilities. Driven by genetic drift and guided by selection, many GRNs that satisfied the same viability function were sorted into parallel lineages, whereas mixing edges between them can lead to fatal pathways and inviable offspring. These results highlight the importance of "functional redundancy" in evolution (*Nowak et al., 1997*; *Láruson et al., 2020*) and agree with earlier studies that suggested alternative regulatory structures can achieve the same phenotype (*True and Haag, 2001*; *Wagner and Wright, 2007*; *Schiffman and Ralph, 2018*). Indeed, both theoretical and empirical studies increasingly support the role of parallel trajectories through fitness landscapes in evolution (*Elmer and Meyer, 2011*; *Bank et al., 2016*; *Ogbunugafor and Eppstein, 2016*; *Langerhans, 2018*).

More importantly, the pathway framework illustrates why degenerate genotypes can reach fix through parallel evolution. Once the alleles are presented as functional pathways connecting an underlying group of proteins, the conjunction between genetic factors and physiological traits is no longer a bipartite mapping; the phenotype, as the collective chemical status of proteins, is a convolution of active signals and external stimuli propagating on the network of genetic pathways. The pathway configuration that satisfies a specific environmental input and phenotypic output is, as a result, not unique. One can thus find numerous, functionally degenerate gene network structures fulfilling the input-output viability relation, as *Figure 4* demonstrates. In addition, taking advantages of basic network analyses, the pathway framework predicts that the number of GRNs generating the same phenotype will increase more than exponentially as the system scales (*Appendix 2*).

The minimal model of GRN evolution we consider encapsulates selection through binary viability, which is essentially a special case of holey adaptive landscapes (*Gavrilets, 1997*). *Gavrilets and Gravner* (*1997*) introduced a multi-locus model where each genotype was independently assigned to one of two fitness levels, whose results suggested that reproductive isolation can arise simply due to the high dimensionality of the genotype space. In a similar vein, our model further connects the high dimensionality of genotypes to complex genetic interactions. Under the pathway framework, inviability originates via the mechanism of hybrid incompatibilities, i.e., allelic combinations that form lethal pathways in a GRN. Furthermore, the pathway framework can be readily extended to include alternative fitness landscapes. For example, *Barton* (*2001*) demonstrated that stabilizing selection can generate reproductive isolation, and the pathway framework can be easily embedded into such a continuous fitness landscape.

Our work supports the latent connection between speciation processes and ancestral genetic variation. Ancient polymorphisms drive genomic divergence and confound inference of evolutionary processes (*Guerrero and Hahn, 2017*). Additionally, these same polymorphisms and the empirical evidence that incompatible alleles often far predate speciation events have recently been consolidated into a "combinatorial" view of speciation (*Marques et al., 2019*). The combinatorial mechanism proposes that, if there was a past admixture event or if standing genetic variation persists, the reassembly of these old genetic variants can facilitate rapid speciation and adaptive radiation. *Marques et al.* found that ancestral genetic variants that had undergone selection–and thus are likely to be beneficial–often have higher allele frequency than *de-novo* mutations. Alternatively, we demonstrate that stochastic loss of accessible pathways resulted in the fixation of incompatible GRNs due to their functional redundancy and high dimensionality. We also observed that the emerging reproductive barriers required the ancestral variation to be greater than a critical amount (*Appendix 3*). Our pathway framework hence adds theoretical support for the role of stable polymorphisms in hybrid incompatibilities, as reviewed in *Cutter* (*2012*). We therefore consider the

456 evolution of regulatory pathways as a parallel mechanism with which ancestral genetic variation
457 can facilitate the appearance of new species.

458     Recent evidence supports our findings that distributed regulatory networks are sources of
459 genetic incompatibilities between closely related taxa. For example, *Morgan et al.* (*2020*) identified
460 a number of disrupted gene expression modules in sub-fertile, hybrid mice and concluded that "hub"
461 genes in these modules played a central role in genetic incompatibility. Additionally, *Rougeux et al.*
462 (*2019*) showed how gene expression was disrupted in hybrids between benthic and limnetic species
463 pairs of Lake Whitefish, *Coregonus clupeaformis* and that genes underlying this disruption were
464 enriched for polymorphisms in the outgroup taxa, the European Whitefish, *Coregonus lavaretus*.
465 This pattern of gene network disruption and standing genetic variation is consistent with our
466 findings from the pathway model. Furthermore, *Guerrero et al.* (*2016*, 2017) found evidence for
467 the role of gene regulatory disruption and the presence of persistent antagonistic interactions in
468 speciation in *Solanum*. Lastly, *Stankowski et al.* (*2019*) found that genetic divergence arose rapidly
469 after population of monkeyflowers were isolated and that the evolution of regulatory-based genetic
470 incompatibilities may have been driven parallel selection pressure from a polymorphic ancestor
471 (*Stankowski et al., 2019*; *Jiggins, 2019*), again the mechanism identified in the pathway framework.

472     Our work is not without important caveats and there are many clear opportunities to advance
473 the pathway framework. First, our model did not include mutation, large-scale genome rearrange-
474 ments, nor whole genome duplication events, which are all known to be important for genetic
475 incompatibles and speciation (*Otto and Whitton, 2000*; *Noor et al., 2001*; *Kirkpatrick and Barton,*
476 *2006*; *Hoffmann and Rieseberg, 2008*; *Guerrero et al., 2012*). Although it is possible to draw some
477 preliminary conclusions regarding the effect of random mutations from our counter-factual sim-
478 ulation that "eliminated" genetic drift, we leave a fuller exploration of mutation for future work.
479 Second, despite the widely documented, asymmetric risk of hybrid breakdown in the heterogametic
480 sex, i.e., Haldanes rule (*Haldane, 1922*; *Coyne and Orr, 1997*; *Delph and Demuth, 2016*), our model
481 considers sexual selection with only a single sex of mating type. Third, both empirical results from
482 yeast (*Bernardes et al., 2017*) and from theoretical, population-genetic models (*Dagilis et al., 2019*)
483 point towards the importance of increased hybrid fitness, i.e., heterosis, even if only temporary,
484 during speciation (*Gavrilets, 2003*). Forth, there are studies that clearly demonstrate the importance
485 of divergent selection in the process of speciation, e.g., (*Nosil et al., 2002*; *Allender et al., 2003*;
486 *Gow et al., 2007*). However, the pathway framework can be readily modified to include divergent
487 selection and will almost certainly result in higher degrees of reproductive isolation. Finally, the
488 relative importance of post-zygotic, genetic incompatibilities in generating and/or maintaining
489 species remains an active area of investigation (*Servedio and Sætre, 2003*; *Rundle and Nosil, 2005*;
490 *Rieseberg and Willis, 2007*; *Magnuson-Ford and Otto, 2012*; *Hopkins, 2013*; *Seehausen et al., 2014*).

491     Our results support a growing body of literature on the theoretical importance of higher-order,
492 genetic interactions in the speciation process (*Johnson and Porter, 2000*; *Palmer and Feldman,*
493 *2009*; *Schiffman and Ralph, 2018*; *Blanckaert et al., 2020*) and are consistent with emerging em-
494 pirical data on genes involved in reproductive isolation (*Seehausen et al., 2014*; *Marques et al.,*
495 *2019*). We support calls for the increased use of high-fidelity simulation models in evolutionary
496 genetics (*Jiggins, 2019*; *Satokangas et al., 2020*), but stress the need for models with interpretable
497 mechanisms and that generate testable hypotheses. For example, our results on the evolution
498 of higher-order incompatibilities could serve as a null model for evaluating empirical data under
499 the relaxed assumption that genes function independently. Only by joining mathematical and
500 computational theory with comparative-level data can we uncover general patterns in speciation
501 and, potentially, resolve long-standing debates in the field.

## Methods

### Numerical Simulations

General Schema and Assumptions

In this work we simulated evolution GRNs in allopatric populations. Throughout evolution, we assumed that individuals had a constant number of loci and thus a fixed number of edges in their GRNs. The underlying set of nodes in GRNs also remained unchanged as we reasoned in Results. We further introduced different categories of nodes/proteins to concrete the space of plausible alleles. Some proteins were presumed to only be present with the environmental stimuli, which were not products of any locus; on the other hand, some other proteins were presumed to have mere physiological effects, and thus they were not capable of activating gene expression. We called them source proteins and target proteins respectively. A plausible allele was therefore labeled by a non-target protein that could activate its expression and a non-source protein that would be synthesized. In our simulations we supposed only one source protein and one target protein.

We considered a naive model of GRN evolution incorporating natural selection, independent assortment and random genetic drift. The environmental condition was set fixed over time and across populations. We assumed that the environment stimulated presence of one protein and it specified another protein with a lethal effect[3]. Viability of individuals was presumably equated to the reciprocal binary state of the lethal protein. Hence given the current generation, individuals were selected such that whoever did not possess a pathway from the environmental stimulus to the lethal protein survived and were able to reproduce.

The survivors then randomly mated and formed the next generation with independent assortment. Here we assumed individuals with haploid-dominant life cycles, where the multicellular haploid stage is evident[4]. Supposed even segregation during meiosis of the diploid zygotes, we modeled the process of independent assortment as follow. Two parental individuals were randomly sampled from the survivors. The set of loci was first randomly partitioned into two groups of equal sizes. The offspring inherited alleles of one group of loci from one of its parents and alleles of the remaining loci from the other parent. Hence half of the edges in the offspring's GRN came from one parent's GRN and the rest was acquired from the other. This procedure was repeated until the next generation had the same constant population size as their predecessors.

Simulations and Parameter Setups

Here we summarize the two different parameter setups in our simulations:

**Setup 1:** We assumed 11 possibly existing proteins in the organism. A generation was composed of 100 individuals with 10 loci each. We generated 100 ancestral populations where individuals' GRNs were randomly sampled from all plausible genotypes. For every ancestral population, we in parallel ran 100 simulations from it, which were regarded as lineages evolving in isolated geo-locations.

**Setup 2:** We assumed 5 possibly existing proteins in the organism. A generation was composed of 16 individuals with 4 loci each. We generated $10^4$ ancestral populations induced from a genetic pool[5] containing all plausible alleles for each locus. For every ancestral population, we in parallel simulated $10^3$ lineages from it.

The randomly generated ancestral populations encapsulate our assumption of ancestral genetic variation, which reflect divergence of gene regulation that has been found in empirical studies (*Gould et al., 2018*). Setup 2 aimed to examine how broadly, in terms of fixed GRNs, evolution can explore in all possibilities. Thus it consisted of a larger amount of simulations starting with unbiased

---

[3]Specifically, they reconciled with the source and the target protein respectively.

[4]During reproduction, specialized haploid cells from two individuals combined and formed a diploid zygote. The zygote experienced meiosis and generated haploid spores, which then developed into multicellular-haploid-stage individuals through mitosis.

[5]We refer a population induced from a genetic pool to a sample among all possible populations that own the same underlying genetic pool.

Manuscript submitted to eLife

546 ancestral populations that were induced from a maximal genetic pool. If not otherwise specified,
547 simulations shown in Results were run under Setup 1.
548    When we inspected reproductive barriers between allopatric populations by interbreeding them,
549 we first sampled 1000 pairs of lineages and then each generated $F_1$ 1000 hybrids. The survival
550 probability of hybrids can then be obtained for all crosses. The same sampling procedure was also
551 applied when we computed the number inter-lineage potential incompatibilities between pairs of
552 allopatric populations.

## Metrics of Reproductive Isolation

554 We introduce a quantitative measure of reproductive isolation between lineages which evolved
555 from a common ancestral population. Given a group of lineages and a chosen pair among them,
556 the reproductive isolation between the pair is defined as the relative difference of hybrid survival

$$I = \frac{p_c - p_h}{p_c} \tag{1}$$

557 where $p_h$ is the survival probability of $F_1$ hybrids, and $p_c$ denotes the average of survival probabilities
558 of all lineages' next generation. A positive value of reproductive isolation $I$ implies that the hybrids
559 have less survivability than the expectation of the offspring. In the extreme case where no hybrid
560 lives, $I = 1$. It therefore serves as an indicator of reproductive barriers between two lineages.
561    Strengths of reproductive barriers among the group of lineages are described through a distribu-
562 tion of reproductive isolation, which can be obtained by sampling pairs of lineages and computing
563 their reproductive isolation $I$. We further introduce two indicators for the existence of reproductive
564 barriers. A quantity named leading reproductive isolation $I^*$ is defined as the 99th percentile of the
565 reproductive isolation distribution. It signals that there is one percent of crosses with reproductive
566 isolation equal or larger than $I^*$. We would also like to raise a caveat that $I^* > 0$ is sufficient for the
567 existence of reproduction barriers but not a necessary condition, due to the possibility of positive
568 $I$ in the distribution even if $I^* \leq 0$. The leading reproductive isolation metric hence summarizes
569 a high level of reproductive barriers that can be found among the lineages. On the other hand,
570 the fraction of positivity in the reproductive isolation distribution serves as a necessity indicator
571 for reproductive barriers, which we denote as $f_p$. The zero-value of $f_p$ implies that none of the
572 crosses generate inviable hybrids more than the anticipation of the offspring and thus the absence
573 of reproductive barriers. Contrarily, a positive $f_p$ does not satisfy existence of barriers considering
574 small reproductive isolation subject to noise. These two indicators are beneficial for us to identify
575 the responsible part of the model to the observed evolutionary consequences.

## Potential Incompatibilities within and between Genetic Pools

577 An intra-lineage incompatibility is a group of alleles in its genetic pool, each of a unique locus, that
578 generates a lethal pathway. In our model those incompatibilities are the only source of inviability,
579 and hence the number of potential incompatibilities provides information about reproductive
580 barriers. Nevertheless, counting the number of potential incompatibilities within a genetic pool
581 through a brute-force manner is computationally intractable. Here we suggest a relatively efficient
582 algorithm when the total number of loci is small. Our strategy is to turn the task into solving a graph
583 problem. The genetic pool can be transformed to an edge-colored network where nodes once more
584 represent possibly existing proteins in the organism. The edges correspond to available alleles
585 in the pool, which are colored by their according loci. A potential incompatibility then becomes
586 a simple path from an environmental input signal to a lethal protein node, with an additional
587 constrain that no edges on the path have the same color. We call such a path an edge-colorful
588 simple path (ECSP).
589    The proposed algorithm, as demonstrated in *Appendix 4* Algorithm 1, counts the number of
590 ECSPs from the source nodes to the targets nodes by having agents propagate on the edge-colored
591 network iteratively. An agents is capable of keeping information of the trajectory, including its

current position on the network, the colors of edges it has traversed and the nodes that it has visited[6]. Initially we deploy one agent on each source node. At every iteration, each agent is substituted by all of its possible successors who are a hop away, such that the hop along with the agent's memory obeys an edge-colorful simple path. Those successors can be deduced from the agent's trajectory information as shown in *Appendix 4* Algorithm 2. The cautiously-designed rule of agent propagation guarantees that the total number of agents locating on the target nodes at the $n$th iteration equals to the number of the desired ECSPs of length $n$. Moreover, since the order of an potential incompatibility is bounded above by the number of genes in the organism, iterations as many as the amount of edge colors in the network are sufficient to obtain a computationally feasible count of all potential incompatibilities. The efficiency of the algorithm can be further improved by, instead of keeping track of numerous agents, monitoring the distribution of agent states over iterations.

The same algorithm can be applied to count the number of inter-lineage incompatibilities as well. In this case the underlying genetic pools of both lineages are transformed into a single edge-colored network, whose edges then consist of alleles in the two pools and are again colored by their according loci. A ECSP on this composite network either only traverses through edges from one of the genetic pools, or it contains alleles from the two different pools. These two scenarios correspond to a incompatibility within and between genetic pools respectively. Therefore, by counting the number of ECSPs on the composite network, and subtracting by the number of potential incompatibilities within the two genetic pools separately, we can compute the number of incompatibilities between the two underlying genetic pools.

## References

**Allender CJ**, Seehausen O, Knight ME, Turner GF, Maclean N. Divergent selection during speciation of Lake Malawi cichlid fishes inferred from parallel radiations in nuptial coloration. Proceedings of the National Academy of Sciences. 2003; 100(24):14074–14079.

**Bank C**, Matuszewski S, Hietpas RT, Jensen JD. On the (un) predictability of a large intragenic fitness landscape. Proceedings of the National Academy of Sciences. 2016; 113(49):14085–14090.

**Barton NH**. The role of hybridization in evolution. Molecular ecology. 2001; 10(3):551–568.

**Bateson W**. Heredity and variation in modern lights. Darwin and modern science. 1909; .

**Bernardes J**, Stelkens R, Greig D. Heterosis in hybrids within and between yeast species. Journal of evolutionary biology. 2017; 30(3):538–548.

**Bikard D**, Patel D, Le Metté C, Giorgi V, Camilleri C, Bennett MJ, Loudet O. Divergent evolution of duplicate genes leads to genetic incompatibilities within A. thaliana. Science. 2009; 323(5914):623–626.

**Blanckaert A**, Bank C, Hermisson J. The limits to parapatric speciation 3: Evolution of strong reproductive isolation in presence of gene flow despite limited ecological differentiation. bioRxiv. 2020; .

**Boyle EA**, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to Omnigenic. Cell. 2017; 169(7):1177 – 1186. http://www.sciencedirect.com/science/article/pii/S0092867417306293, doi: https://doi.org/10.1016/j.cell.2017.05.038.

**Brideau NJ**, Flores HA, Wang J, Maheshwari S, Wang X, Barbash DA. Two Dobzhansky-Muller genes interact to cause hybrid lethality in Drosophila. science. 2006; 314(5803):1292–1295.

**Castillo DM**, Barbash DA. Moving speciation genetics forward: modern techniques build on foundational studies in Drosophila. Genetics. 2017; 207(3):825–842.

**Chae E**, Bomblies K, Kim ST, Karelina D, Zaidem M, Ossowski S, Martín-Pizarro C, Laitinen RE, Rowan B, Tenenboim H, Lechner S, Demar M, Habring-Müller A, Lanz C, Rätsch G, Weigel D. Species-wide Genetic Incompatibility Analysis Identifies Immune Genes as Hot Spots of Deleterious Epistasis. Cell. 2014; 159(6):1341 – 1351. http://www.sciencedirect.com/science/article/pii/S0092867414013762, doi: https://doi.org/10.1016/j.cell.2014.10.049.

---

[6]In Algorithm 1, the NEW-AGENT procedure creates an agent instance given its position, visited colors and nodes accordingly. This trajectory information is also accessible fields of the agent instance.

639  **Coyne JA**, Allen Orr H. The evolutionary genetics of speciation. Philosophical Transactions of the Royal Society
640      of London Series B: Biological Sciences. 1998; 353(1366):287–305.

641  **Coyne JA**, Orr HA. "Patterns of speciation in Drosophila" revisited. Evolution. 1997; 51(1):295–303.

642  **Cutter AD**. The polymorphic prelude to Bateson–Dobzhansky–Muller incompatibilities. Trends in ecology &
643      evolution. 2012; 27(4):209–218.

644  **Dagilis AJ**, Kirkpatrick M, Bolnick DI. The evolution of hybrid fitness during speciation. PLoS genetics. 2019;
645      15(5).

646  **Dagilis AJ**, Matute DR. Incompatibilities between emerging species. Science. 2020; 368(6492):710–711.

647  **Davidich M**, Bornholdt S. The transition from differential equations to Boolean networks: A case study in
648      simplifying a regulatory network model. Journal of Theoretical Biology. 2008; 255(3):269 – 277. http://www.
649      sciencedirect.com/science/article/pii/S0022519308003652, doi: https://doi.org/10.1016/j.jtbi.2008.07.020.

650  **Davies B**, Hatton E, Altemose N, Hussin JG, Pratto F, Zhang G, Hinch AG, Moralli D, Biggs D, Diaz R, et al.
651      Re-engineering the zinc fingers of PRDM9 reverses hybrid sterility in mice. Nature. 2016; 530(7589):171–176.

652  **De Bruijn NG**. A combinatorial problem. In: *Proc. Koninklijke Nederlandse Academie van Wetenschappen*, vol. 49;
653      1946. p. 758–764.

654  **Delph LF**, Demuth JP. Haldane's rule: genetic bases and their empirical support. Journal of Heredity. 2016;
655      107(5):383–391.

656  **Dobzhansky T**. Studies on hybrid sterility. II. Localization of sterility factors in Drosophila pseudoobscura
657      hybrids. Genetics. 1936; 21(2):113.

658  **Duranton M**, Allal F, Valière S, Bouchez O, Bonhomme F, Gagnaire PA. The contribution of ancient admixture to
659      reproductive isolation between European sea bass lineages. BioRxiv. 2019; p. 641829.

660  **Elmer KR**, Meyer A. Adaptation in the age of ecological genomics: insights from parallelism and convergence.
661      Trends in ecology & evolution. 2011; 26(6):298–306.

662  **Gavrilets S**. Evolution and speciation on holey adaptive landscapes. Trends in ecology & evolution. 1997;
663      12(8):307–312.

664  **Gavrilets S**. Perspective: models of speciation: what have we learned in 40 years? Evolution. 2003; 57(10):2197–
665      2215.

666  **Gavrilets S**, Gravner J. Percolation on the fitness hypercube and the evolution of reproductive isolation. Journal
667      of theoretical biology. 1997; 184(1):51–64.

668  **Gould BA**, Chen Y, Lowry DB. Gene regulatory divergence between locally adapted ecotypes in their native
669      habitats. Molecular Ecology. 2018; 0(0). https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.14852, doi:
670      10.1111/mec.14852.

671  **Gow JL**, Peichel CL, Taylor EB. Ecological selection against hybrids in natural populations of sympatric threespine
672      sticklebacks. Journal of evolutionary biology. 2007; 20(6):2173–2180.

673  **Guerrero RF**, Hahn MW. Speciation as a sieve for ancestral polymorphism. Molecular Ecology. 2017; 26(20):5362–
674      5368.

675  **Guerrero RF**, Muir CD, Josway S, Moyle LC. Pervasive antagonistic interactions among hybrid incompatibility
676      loci. PLoS genetics. 2017; 13(6):e1006817.

677  **Guerrero RF**, Posto AL, Moyle LC, Hahn MW. Genome-wide patterns of regulatory divergence revealed by
678      introgression lines. Evolution. 2016; 70(3):696–706.

679  **Guerrero RF**, Rousset F, Kirkpatrick M. Coalescent patterns for chromosomal inversions in divergent populations.
680      Philosophical Transactions of the Royal Society B: Biological Sciences. 2012; 367(1587):430–438.

681  **Haldane J**. Sex ratio and unisexual sterility in hybrid animals. Journal of genetics. 1922; 12(2):101–109.

682  **Han F**, Lamichhaney S, Grant BR, Grant PR, Andersson L, Webster MT. Gene flow, ancient polymorphism, and
683      ecological adaptation shape the genomic landscape of divergence among Darwin's finches. Genome research.
684      2017; 27(6):1004–1015.

Manuscript Submitted to eLife

**Hoffmann AA**, Rieseberg LH. Revisiting the impact of inversions in evolution: from population genetic markers to drivers of adaptive shifts and speciation? Annual review of ecology, evolution, and systematics. 2008; 39:21–42.

**Hopkins R**. Reinforcement in plants. New Phytologist. 2013; 197(4):1095–1103.

**Hopkins R**, Rausher MD. Identification of two genes causing reinforcement in the Texas wildflower Phlox drummondii. Nature. 2011; 469(7330):411–414.

**Jamie GA**, Meier JI. The Persistence of Polymorphisms across Species Radiations. Trends in Ecology & Evolution. 2020; .

**Jiggins CD**. Can genomics shed light on the origin of species? PLoS biology. 2019; 17(8):e3000394.

**Johnson NA**, Porter AH. Rapid speciation via parallel, directional selection on regulatory genetic pathways. Journal of Theoretical Biology. 2000; 205(4):527–542.

**Kaeuffer R**, Peichel CL, Bolnick DI, Hendry AP. Parallel and nonparallel aspects of ecological, phenotypic, and genetic divergence across replicate population pairs of lake and stream stickleback. Evolution: International Journal of Organic Evolution. 2012; 66(2):402–418.

**Kalirad A**, Azevedo RBR. Spiraling Complexity: A Test of the Snowball Effect in a Computational Model of RNA Folding. Genetics. 2017; 206(1):377–388. http://www.genetics.org/content/206/1/377, doi: 10.1534/genetics.116.196030.

**Kelly DE**, Hansen ME, Tishkoff SA. Global variation in gene expression and the value of diverse sampling. Current opinion in systems biology. 2017; 1:102–108.

**Kirkpatrick M**, Barton N. Chromosome inversions, local adaptation and speciation. Genetics. 2006; 173(1):419–434.

**Kuzmin E**, VanderSluis B, Wang W, Tan G, Deshpande R, Chen Y, Usaj M, Balint A, Mattiazzi Usaj M, van Leeuwen J, Koch EN, Pons C, Dagilis AJ, Pryszlak M, Wang JZY, Hanchard J, Riggi M, Xu K, Heydari H, San Luis BJ, et al. Systematic analysis of complex genetic interactions. Science. 2018; 360(6386). http://science.sciencemag.org/content/360/6386/eaao1729, doi: 10.1126/science.aao1729.

**Langerhans RB**. Predictability and parallelism of multitrait adaptation. Journal of Heredity. 2018; 109(1):59–70.

**Langfelder P**, Horvath S. WGCNA: an R package for weighted correlation network analysis. BMC bioinformatics. 2008; 9(1):559.

**Láruson ÁJ**, Yeaman S, Lotterhos KE. The Importance of Genetic Redundancy in Evolution. Trends in Ecology & Evolution. 2020; .

**Livingstone K**, Olofsson P, Cochran G, Dagilis A, MacPherson K, Seitz KA. A stochastic model for the development of Bateson–Dobzhansky–Muller incompatibilities that incorporates protein interaction networks. Mathematical Biosciences. 2012; 238(1):49 – 53. http://www.sciencedirect.com/science/article/pii/S0025556412000491, doi: https://doi.org/10.1016/j.mbs.2012.03.006.

**Lowry DB**, Hernandez K, Taylor SH, Meyer E, Logan TL, Barry KW, Chapman JA, Rokhsar DS, Schmutz J, Juenger TE. The genetics of divergence and reproductive isolation between ecotypes of Panicum hallii. New Phytologist. 2015; 205(1):402–414.

**Lowry DB**, Modliszewski JL, Wright KM, Wu CA, Willis JH. The strength and genetic basis of reproductive isolating barriers in flowering plants. Philosophical Transactions of the Royal Society B: Biological Sciences. 2008; 363(1506):3009–3021.

**Magnuson-Ford K**, Otto SP. Linking the investigations of character evolution and species diversification. The American Naturalist. 2012; 180(2):225–245.

**Marques DA**, Meier JI, Seehausen O. A combinatorial view on speciation and adaptive radiation. Trends in ecology & evolution. 2019; .

**Meier JI**, Marques DA, Mwaiko S, Wagner CE, Excoffier L, Seehausen O. Ancient hybridization fuels rapid cichlid fish adaptive radiations. Nature communications. 2017; 8:14363.

**Mogil LS**, Andaleon A, Badalamenti A, Dickinson SP, Guo X, Rotter JI, Johnson WC, Im HK, Liu Y, Wheeler HE. Genetic architecture of gene expression traits across diverse populations. PLoS genetics. 2018; 14(8):e1007586.

**Morgan K**, Harr B, White MA, Payseur BA, Turner LM. Disrupted gene networks in subfertile hybrid house mice. Molecular biology and evolution. 2020; 37(6):1547–1562.

**Moyle LC**, Nakazato T. Comparative genetics of hybrid incompatibility: sterility in two Solanum species crosses. Genetics. 2008; 179(3):1437–1453.

**Muller H**. Isolating mechanisms, evolution, and temperature. In: *Biol. Symp.*, vol. 6; 1942. p. 71–125.

**Nelson TC**, Cresko WA. Ancient genomic variation underlies repeated ecological adaptation in young stickleback populations. Evolution Letters. 2018; 2(1):9–21.

**Noor MA**, Feder JL. Speciation genetics: evolving approaches. Nature Reviews Genetics. 2006; 7(11):851–861.

**Noor MA**, Grams KL, Bertucci LA, Reiland J. Chromosomal inversions and the reproductive isolation of species. Proceedings of the National Academy of Sciences. 2001; 98(21):12084–12088.

**Nosil P**, Crespi BJ, Sandoval CP. Host-plant adaptation drives the parallel evolution of reproductive isolation. Nature. 2002; 417(6887):440–443.

**Nosil P**, Schluter D. The genes underlying the process of speciation. Trends in ecology & evolution. 2011; 26(4):160–167.

**Nowak MA**, Boerlijst MC, Cooke J, Smith JM. Evolution of genetic redundancy. Nature. 1997; 388(6638):167–171.

**Ogbunugafor CB**, Eppstein MJ. Competition along trajectories governs adaptation rates towards antimicrobial resistance. Nature ecology & evolution. 2016; 1(1):1–8.

**Orr HA**. The population genetics of speciation: the evolution of hybrid incompatibilities. Genetics. 1995; 139(4):1805–1813. http://www.genetics.org/content/139/4/1805.

**Otto SP**, Whitton J. Polyploid incidence and evolution. Annual review of genetics. 2000; 34(1):401–437.

**Palmer ME**, Feldman MW. DYNAMICS OF HYBRID INCOMPATIBILITY IN GENE NETWORKS IN A CONSTANT ENVIRONMENT. Evolution. 2009; 63(2):418–431. https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1558-5646.2008.00577.x, doi: 10.1111/j.1558-5646.2008.00577.x.

**Powell DL**, García-Olazábal M, Keegan M, Reilly P, Du K, Díaz-Loyo AP, Banerjee S, Blakkan D, Reich D, Andolfatto P, et al. Natural hybridization reveals incompatible alleles that cause melanoma in swordtail fish. Science. 2020; 368(6492):731–736.

**Presgraves DC**. The molecular evolutionary basis of species formation. Nature Reviews Genetics. 2010; 11(3):175–180.

**Rieseberg LH**, Sinervo B, Linder CR, Ungerer MC, Arias DM. Role of gene interactions in hybrid speciation: evidence from ancient and experimental hybrids. Science. 1996; 272(5262):741–745.

**Rieseberg LH**, Willis JH. Plant speciation. science. 2007; 317(5840):910–914.

**Rougeux C**, Gagnaire PA, Praebel K, Seehausen O, Bernatchez L. Polygenic selection drives the evolution of convergent transcriptomic landscapes across continents within a Nearctic sister species complex. Molecular ecology. 2019; 28(19):4388–4403.

**Rundle HD**, Nosil P. Ecological speciation. Ecology letters. 2005; 8(3):336–352.

**Ryu KH**, Huang L, Kang HM, Schiefelbein J. Single-cell RNA sequencing resolves molecular relationships among individual plant cells. Plant physiology. 2019; 179(4):1444–1456.

**Satokangas I**, Martin S, Helanterä H, Saramäki J, Kulmuni J. Multi-locus interactions and the build-up of reproductive isolation. arXiv preprint arXiv:200513790. 2020; .

**Schiffman JS**, Ralph PL. System drift and speciation. bioRxiv. 2018; https://www.biorxiv.org/content/early/2018/01/26/231209, doi: 10.1101/231209.

**Schlitt T**, Brazma A. Current approaches to gene regulatory network modelling. BMC bioinformatics. 2007; 8(S6):S9.

**Schluter D**. Evidence for Ecological Speciation and Its Alternative. Science. 2009; 323(5915):737–741. http://science.sciencemag.org/content/323/5915/737, doi: 10.1126/science.1160006.

Manuscript submitted to eLife

**Seehausen O**, Butlin RK, Keller I, Wagner CE, Boughman JW, Hohenlohe PA, Peichel CL, Saetre GP, Bank C, Brännström Å, et al. Genomics and the origin of species. Nature Reviews Genetics. 2014; 15(3):176–192.

**Servedio MR**, Sætre GP. Speciation as a positive feedback loop between postzygotic and prezygotic barriers to gene flow. Proceedings of the Royal Society of London Series B: Biological Sciences. 2003; 270(1523):1473–1479.

**Sicard A**, Kappel C, Josephs EB, Lee YW, Marona C, Stinchcombe JR, Wright SI, Lenhard M. Divergent sorting of a balanced ancestral polymorphism underlies the establishment of gene-flow barriers in Capsella. Nature communications. 2015; 6:7960.

**Stankowski S**, Chase MA, Fuiten AM, Rodrigues MF, Ralph PL, Streisfeld MA. Widespread selection and gene flow shape the genomic landscape during a radiation of monkeyflowers. PLoS biology. 2019; 17(7):e3000391.

**Tong AHY**, Lesage G, Bader GD, Ding H, Xu H, Xin X, Young J, Berriz GF, Brost RL, Chang M, et al. Global mapping of the yeast genetic interaction network. science. 2004; 303(5659):808–813.

**True JR**, Haag ES. Developmental system drift and flexibility in evolutionary trajectories. Evolution & development. 2001; 3(2):109–119.

**Turner LM**, White MA, Tautz D, Payseur BA. Genomic Networks of Hybrid Sterility. PLOS Genetics. 2014 02; 10(2):1–23. https://doi.org/10.1371/journal.pgen.1004162, doi: 10.1371/journal.pgen.1004162.

**Tyler AL**, Ji B, Gatti DM, Munger SC, Churchill GA, Svenson KL, Carter GW. Epistatic networks jointly influence phenotypes related to metabolic disease and gene expression in diversity outbred mice. Genetics. 2017; 206(2):621–639.

**Vaid N**, Laitinen RA. Diverse paths to hybrid incompatibility in Arabidopsis. The Plant Journal. 2019; 97(1):199–213.

**Wagner A**, Wright J. Alternative routes and mutational robustness in complex regulatory networks. Biosystems. 2007; 88(1-2):163–172.

**Wang B**, Mojica JP, Perera N, Lee CR, Lovell JT, Sharma A, Adam C, Lipzen A, Barry K, Rokhsar DS, et al. Ancient polymorphisms contribute to genome-wide variation by long-term balancing selection and divergent sorting in Boechera stricta. Genome biology. 2019; 20(1):126.

**Wang RJ**, White MA, Payseur BA. The pace of hybrid incompatibility evolution in house mice. Genetics. 2015; 201(1):229–242.

**Wilf HS**. Generating functionology. Elsevier; 2013.

**Wittbrodt J**, Adam D, Malitschek B, Mäueler W, Raulf F, Telling A, Robertson SM, Schartl M. Novel putative receptor tyrosine kinase encoded by the melanoma-inducing Tu locus in Xiphophorus. Nature. 1989; 341(6241):415–421.

**Wolf JB**, Lindell J, Backström N. Speciation genetics: current status and evolving approaches. Phil Trans R Soc B. 2010; 365:1717—-1733.

**Yamamoto E**, Takashi T, Morinaka Y, Lin S, Wu J, Matsumoto T, Kitano H, Matsuoka M, Ashikari M. Gain of deleterious function causes an autoimmune response and Bateson–Dobzhansky–Muller incompatibility in rice. Molecular Genetics and Genomics. 2010; 283(4):305–315.

Manuscript submitted to eLife

## Appendix 1

### Hybrid inviability against a single incompatibility

Here we analytically evaluate the probability that a hybrid is inviable presuming that multiple incompatibilities are rarely embedded in two parental gene regulatory networks. In addition, this naive analysis explains the pattern of RI distribution, *Figure 5a* in the main text.

Assume that there is only on incompatibility $\mathcal{I}$ between the two parental gene networks $G_1$ and $G_2$. For convenience we suppose there are an even number of loci in the organisms, denoted by $2m$, and let the incompatibility $\mathcal{I}$ be of order $k-1$ so it consists of $k$ alleles to form a lethal combination. We also suppose that, among the $k$ alleles in $\mathcal{I}$, $k_1$ of them come from $G_1$ and the other $k_2$ alleles are from $G_2$.

Following the rule of recombination between haploid GRNs in our model, the hybrid is generated by randomly segregating alleles of $m$ loci from $G_1$ and then mixing with alleles of the other $m$ loci from $G_2$. Hence if $m < k_1$ or $m < k_2$, then there is no chance that the incompatibility $\mathcal{I}$ appears in the hybrid. Otherwise, among all plausible segregation, we can compute the number of achievable ways that the $k_1$ and $k_2$ alleles from $G_1$ and $G_2$ respectively are sorted into the hybrid. The probability that the hybrid is inviable due to the only incompatibility $\mathcal{I}$ is thus
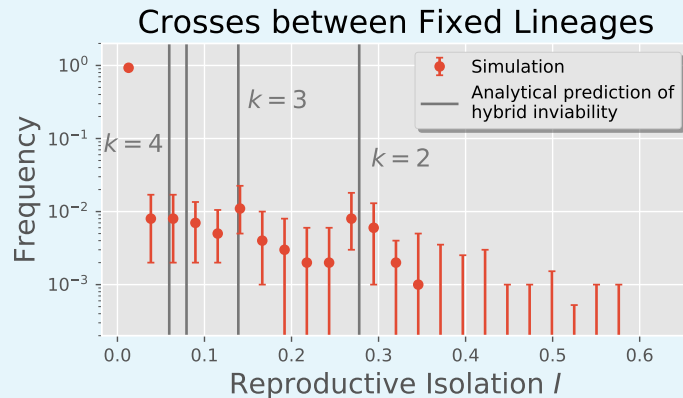
$$P(\mathcal{I}) = \begin{cases} \dfrac{\binom{2m-k}{m-k_1}}{\binom{2m}{m}}, & \text{if } k_1, k_2 \leq m \\ 0, & \text{otherwise} \end{cases} \tag{2}$$

If we further assume that $m \gg 1$ and $m \gg k$, applying the Stirling's approximation we have an estimate of the hybrid inviability

$$P(\mathcal{I}) = \frac{m!m!(2m-k)!}{(m-k_1)!(m-k_2)!(2m)!} \approx 2^{-k} \tag{3}$$

This plain derivation shows that, should there be only one incompatibility concealing between two parental GRNs, the survivability of a hybrid is predominantly determined by the order of the incompatibility.

Here *Figure 1* shows good agreement between our analytical prediction of hybrid inviability and the "bulges" from the observed RI distribution. Our simple derivation explains the higher likelihood of certain RI levels relative to their neighboring regions. It also manifests how the discreteness nature of hybrid incompatibilities shapes the RI distribution and that this characteristic has major effects on the strength of reproductive barriers.



**Appendix 1 Figure 1.** Comparison between the uncovered RI distribution in our simulations and the predicted hybrid inviability *Equation 2*.

## Appendix 2

### Estimating functional redundancy of GRNs under extreme selection

Our pathway framework not only resonates with existing studies of the functional redundancy of GRNs (*True and Haag, 2001*; *Wagner and Wright, 2007*; *Schiffman and Ralph, 2018*; *Láruson et al., 2020*), but it also estimates how many GRNs generate a given phenotype under the Boolean-state assumption. Here we consider an extreme case where every protein is either required present or absent, except those that are stimulated by the environment. This scenario depicts a strong selection force, and a weaker selection can be easily reached by relaxing the phenotypic constraint on proteins. Note that this extreme scenario hence provides a lower bound of the number of GRNs that produce the same phenotype.

Suppose there are $n_+$ and $n_-$ proteins that are required present and absent respectively, and let there be $n_0$ present-state proteins due to the environmental stimuli. A GRN that generates this given phenotype can be viewed as a composition of two parts: First, it contains alleles, i.e., edges, building up pathways from any of the $n_0$ stimulated proteins to every of the $n_+$ required-present proteins. Second, edges associated with the required-absent proteins, if any, must not be alleles activated by the required-present/stimulated proteins and producing the required-absent ones. Assuming $m$ haploid loci ($m \geq n_+$), the number of GRNs generating the given phenotype is

$$f(m, n_0, n_+, n_-) = \sum_{k=n_+}^{m} \binom{m}{k} \left[ n_- \left( n_- + n_+ \right) \right]^{m-k} f(k, n_0, n_+, 0), \tag{4}$$

where $f(k, n_0, n_+, 0)$ corresponds to the special case where no protein is required absent, it is equivalently the number of directed, edge-labeled graphs with $n_0 + n_+$ nodes and $k$ edges such that every of the $n_+$ nodes are reachable from any of the $n_0$ nodes.

Although one may compute compute $f(k, n_0, n_+, 0)$ through a recursive relation generalized from existing literature (e.g., *Wilf, 2013*), an analytical solution is hardly accessible. Here we instead assess the lower and upper bound of $f(m, n_0, n_+, n_-)$. First, $f(m, n_0, n_+, n_-)$ accounts for all graphs satisfying the reachability criterion, and it is bounded below by the amount of those graphs which are also forests[a]. Finding all such forests is equivalent to finding all possibilities to grow a network from $n_0$ initial nodes, where edges are added incrementally, pointing from an existing node to a not-yet-existing (newly added) one. So we have

$$f(m, n_0, n_+, n_-) \geq \binom{m}{n_+} \left[ n_- \left( n_- + n_+ \right) \right]^{m-n_+} f(n_+, n_0, n_+, 0)$$

$$= \binom{m}{n_+} \left[ n_- \left( n_- + n_+ \right) \right]^{m-n_+} \frac{n_+! \, (n_0 + n_+)!}{n_0!}. \tag{5}$$

Second, incrementally adding $k - n_+$ edges to every of such forests is essentially an enumeration of the $k$-edge graphs. This process generates all possible graphs with $k$ labeled-edges, but they might be over-counted since adding edges to two different forests can produce the same graph. Computing all possible ways to add $k$ edges to every forest satisfying the reachability criterion leads to an upper bound of $f(k, n_0, n_+, 0)$:

$$f(k, n_0, n_+, 0) \leq \left[ n_+ \left( n_0 + n_+ \right) \right]^{k-n_+} \frac{n_+! \, (n_0 + n_+)!}{n_0!}, \tag{6}$$

and $f(m, n_0, n_+, n_-)$ is hence bounded above by

$$f(m, n_0, n_+, n_-) \leq \sum_{k=n_+}^{m} \binom{m}{k} \left[ n_- \left( n_- + n_+ \right) \right]^{m-k} \left[ n_+ \left( n_0 + n_+ \right) \right]^{k-n_+} \frac{n_+! \, (n_0 + n_+)!}{n_0!}. \tag{7}$$

898  Combined, under an extreme scenario where the binary states of all proteins are con-
899  strained, we see super-linearly or even exponentially many GRNs generating the same
900  phenotype. The pathway framework therefore concludes that, for any phenotype derived
901  from the binary states of proteins, the number of functionally redundant GRNs grows faster
902  than super-linearly/exponentially as the system scales.

---

[a]A *forest* is a graph which only has trees as its connected components, where trees are graphs without cycles.
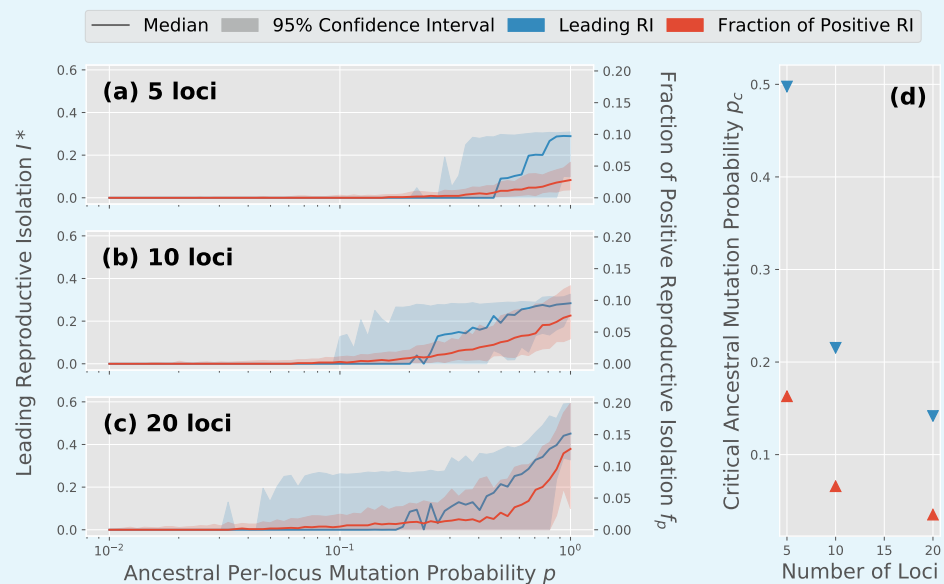
Manuscript submitted to eLife

### Appendix 3

#### Reproductive barriers and ancestral genetic variation

Here we demonstrate our examination on how the extent of ancestral genetic variation influences the appearance and strength of reproductive barriers. To begin with, we designed a pipeline to produce ancestral populations whose amount of genetic variation are tunable. A fixed population was first obtained from our GRN evolution model starting with randomly generated individual GRNs. For every locus, the allele might then mutate into any other possible allele with a per-locus mutation probability $p$. The resulting population was regarded as the ancestral population, where the mutation probability $p$ became a tunable parameter to assess the degree of ancestral variation.

We followed the same methodology to simulate generational dynamics of GRNs and to compute reproductive isolation between allopatric lineages as in the main text. *Figure 1a-c* below shows, for different number of loci, the reproductive barriers consequent to the varying ancestral mutation probability $p$. Here we present two indicators of barriers: the leading RI (blue, left axis) and the fraction of positive RI (red, right axis). On a first glance the simulations evince that, for a organism with a larger number of loci, the barriers only required a smaller ancestral mutation probability yet more apparent barriers were observed.

*Figure 1a-c* furthermore suggest some critical level of ancestral variation associated with the constant population size, such that reproductive barriers would hardly appear between lineages evolving from an ancestral population with less polymorphisms. We quantify the critical level of genetic variation through a critical mutation probability $p_c$; this is the smallest ancestral mutation probability with which a barrier indicator has non-zero median value. Nevertheless, due to the lack of a both sufficient and necessary indicator, we could only estimate the interval that this critical level fell into. The critical level of ancestral variation would be bounded above by $p_c$ for the leading RI (a sufficient indicator of barriers) and bounded below by one for the fraction of positive RI (a necessary indicator of barriers). *Figure 1d* presents the interval estimation that the critical ancestral variation fell into for organisms with different number of loci.



**Appendix 3 Figure 1.** Varying the extent ancestral variation and its corresponding strength of reproductive barriers. The GRN evolution was simulated under Setup 1 described in Methods. **(a-c)** Indicators of barriers for 5, 10 and 20 loci. **(d)** Estimation of their critical level of ancestral variation.

## Appendix 4

### Algorithms of counting potential incompatibilities

---
**Algorithm 1** COUNT-ECSP
---

**Require:** A set of source nodes $S$; a set of target nodes $T$; a map $I$ from nodes to their incident outgoing edges; a set of path lengths of interests $L$.

**Ensure:** A map $C$ from $L$ to the number of edge-colorful simple paths from $S$ to $T$, which are of the corresponding length.

1: $C \leftarrow$ an empty map
2: $l_{max} \leftarrow$ the largest element of $L$
3: $A \leftarrow$ an empty list                                     ▷ Initialize agents.
4: **for all** node $s \in S$ **do** $A$.INSERT(NEW-AGENT($s$, Ø, $\{s\}$))
5: **end for**
6: **for** $l \leftarrow 1$ to $l_{max}$ **do**      ▷ Iterate over the number of hops agents have made from the source nodes.
7:     $n \leftarrow 0$
8:     $N \leftarrow$ an empty list                               ▷ Update the list of agents.
9:     **for all** agent $a \in A$ **do**
10:         **for all** agent $a' \in$ NEXT-POSSIBILITIES($a$, $I$) **do**
11:             $N$.INSERT($a'$)
12:             **if** $a'.position \in T$ **then** $n \leftarrow n + 1$
13:             **end if**
14:         **end for**
15:     **end for**
16:     $A \leftarrow N$
17:     **if** $l \in L$ **then** $C$.INSERT($l$, $n$)                    ▷ Update counting.
18:     **end if**
19: **end for**
20: **return** $C$

---
**Algorithm 2** NEXT-POSSIBILITIES
---

**Require:** An agent $a$; a map $I$ from nodes to their incident outgoing edges.

**Ensure:** A set $P$ of agents who are of all the possible states that can be reached through a hop from the given agent $a$, such that

1. The hop only goes through an edge of a color that has not been visited by the agent.
2. The position after the hop has not been visited by the agent.

1: $P \leftarrow$ an empty set
2: **for all** edge $e \in I$.GET($a$) **do**
3:     **if** $e.color \notin a.colors\text{-}visited$ and $e.target \notin a.nodes\text{-}visited$ **then**
4:         $a' \leftarrow$ NEW-AGENT($e.target$, $a.colors\text{-}visited \cup \{e.color\}$, $a.nodes\text{-}visited \cup \{e.target\}$)
5:         $P$.INSERT($a'$)
6:     **end if**
7: **end for**
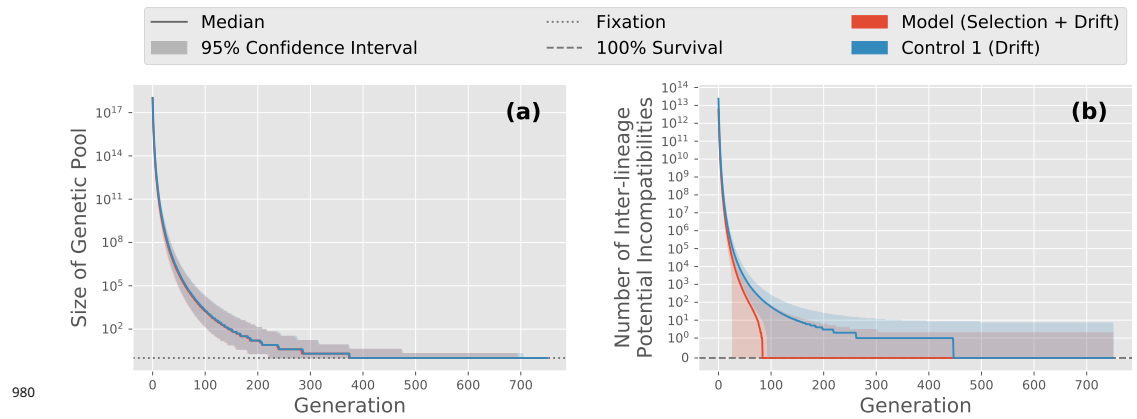8: **return** $P$

Manuscript submitted to eLife



980

**Figure 7–Figure supplement 1. (a)** The size the underlying genetic pool continually shrank until there was only one accessible genotype. At this stage a population fixated a single GRN, and no significant difference was found between the model and the control scenario without selection, i.e., drift only. **(b)** In our model, inter-lineage incompatibilities persisted throughout evolution (red), which accounts for the sustained confidence interval of their abundance even after populations reach fixation. Interestingly, in the control scenario where natural selection was silenced, inter-lineage incompatibilities were eliminated at a slower pace. We hypothesize that due to the lack of guidance by selection, inter-lineage incompatibilities only became inaccessible through random genetic drift. This scenario led to fatal allelic combinations that were more persistent than those in the model and hence stronger reproductive barriers were observed.