

Architecture of the Genotype-Phenotype Map and the Coevolution of Complexity

Bhaskar Kumawat¹ and Luis Zaman^{2,3}

¹Undergraduate Program, Indian Institute of Science, Bangalore, India - 560012

²Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109, USA

³Center for the Study of Complex Systems, University of Michigan, Ann Arbor, MI 48109, USA
{kumawatb, zamanlh}@umich.edu

Abstract

The addition of parasites to a host population can drive an escalation in the host population’s phenotypic complexity – even in the absence of a direct fitness advantage for this increase. Parasites restrict certain regions of the genotype space, decreasing the fitness and the probability of survival of particular host phenotypes. While many artificial life frameworks model a direct correlation between genotype and fitness, the structure of genotype-phenotype maps can have important effects on evolutionary dynamics. Using a simple coarse-grained model for phenotypic transitions during evolution, we show that the escalation in phenotypic complexity under neutral co-evolution is dependent on the structure of the genotype-phenotype map. We discuss these results using the metaphor of evolutionary spandrels and highlight how these structural considerations might allow us to capture biological phenomena more accurately.

Introduction

Extant life on earth uses a common method of information storage and transfer – the sequence of DNA bases comprising the genome. This information is converted into mediators of chemical and physical interactions in the cell – proteins and RNAs – which in turn determine the phenotype of an organism embedded within a particular environment. A primary goal of artificial life is to recreate processes that mimic living systems. At the same time, predicting the outcome of such processes in biological systems, but using artificial life frameworks, has been a fruitful endeavor (Goldsby et al., 2014; Nelson and Sanford, 2011; Wilke et al., 2001). Many digital frameworks model genetic information at a level where it is directly involved in phenotype determination. However, in a cell for example, the phenotype is an emergent property of many interacting processes at multiple scales. How structural properties of genotype-phenotype (GP) maps affect the processes and outcomes of evolution deserves much more attention. In this work, we will investigate how architectural features of the GP-map influences the coevolution of complexity using a simple theoretical model.

Specifically, in neutral environments where all phenotypes are equally fit, we find that a population of organisms will distribute evenly over the genotype space. Phenotypes

encoded by a larger number of genotypes will be more likely to appear solely due to the unbiased nature of the mutation process. Addition of a coevolving parasite population, however, throws this equilibrated host population into disarray. The host population must now switch between phenotypes to escape infection and therefore survive. Indeed, these co-evolutionary dynamics have been harnessed to improve evolutionary algorithms many times (Hillis, 1990; Floreano and Nolfi, 1997; Miikkulainen and Stanley, 2004; Wagner et al., 2020; Watson and Pollack, 2001). Envisioning this “arms race” over a limited number of states, we hypothesize and show that the structure of the mapping between genotypes and phenotypes plays an important role in the coevolutionary escalation in complexity (and diversity) and, therefore, might be the ultimate evolutionary spandrel.

Model

Coarse-graining the genotype network

Consider a genotype space consisting of bit-strings of length L . This space consists of a total of 2^L sequences which can be divided into different phenotypes depending on how the chemistry in the system allows this information to be decoded into physical traits (phenotypes) of the organism. This process, which we will hereby refer to as *partitioning* of the genotype space, can be performed in different ways. On top of this set-like structure of the partitioned genotype space, there is a secondary network-like structure that arises out of mutational proximity between the genotypes. Connecting these partitioned genotypes through single-mutation edges constructs the partitioned genotype network (Figure 1). The mutation process is, in general, unbiased across the entire space and the probability of traversing a particular path through the network in a single generation can be given in terms of the per-site mutation rate (m) and the mutational distance (l) between any two given genotypes.

$$p_{g_j \rightarrow g_i} = (1 - m)^{L-l(g_j, g_i)} m^{l(g_j, g_i)} \quad (1)$$

Using the values obtained from the above equation, we can form a genotype-transition matrix (G) consisting of

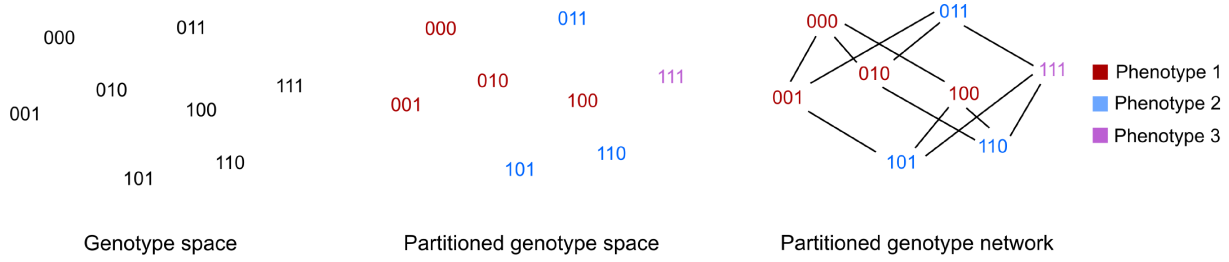


Figure 1: Partitioning of 3-bit genotypes into phenotypes and the network structure of the genotype space. Considering all binary strings (bit-strings) of a particular length L (here, $L=3$), we identify the space of all genotypes (left). Next, we partition these genotypes into phenotypes (here, $N_P=3$) depicted as different colors (center). Finally, we connect bit-strings accessible by a single mutation with an edge to construct the complete partitioned genotype network (right).

transition probabilities between genotypes through mutation.

G_{ij} = probability of transition from genotype g_j to g_i in a generation

This matrix quantifiably captures the network-like structure of the genotype space. The partitioning into phenotypes can be described by a separate matrix, the partitioning-scheme matrix (S), as follows.

$$S_{ij} = \begin{cases} 1 & , \text{ if genotype } g_j \text{ belongs to phenotype } i \\ 0 & , \text{ otherwise} \end{cases} \quad (2)$$

Assuming that the mutation and selection processes do not discriminate between genotypes under the same phenotype and that all genotype transitions are equally probable, we can coarse-grain the genotype transition matrix into a phenotype-transition matrix (P).

$$P = SGS^T \quad (3)$$

Where S^T is the transpose of the partitioning-scheme matrix and the matrix P is normalized column-wise. Each term of the matrix P approximates the probability of transition from one phenotype to another in a single generation. While using genotype transitions is not feasible due to the sheer size of the genotype space for non-trivial values of L , calculation of the phenotype-transition matrix makes computational simulation of host and parasite populations under these phenotypes possible. We have verified that the coarse-grained version gives similar results as the fine-grained genotype-by-genotype model under the same conditions.

Phenotypic Complexity

Next, we turn to the problem of defining phenotypic complexity in an operational and meaningful way. Because complexity is a term used to cover a set of related concepts, a singular unambiguous definition is elusive. We are interested in capturing a notion of evolutionary difficulty, or in

other words, the amount of surprise upon observing a trait that did in fact evolve (Wagner, 2017)

Certain traits are easier to innovate and can thus be specified flexibly by a large number of similar genotypes. On the other hand, some traits might require a very specific genotype sequence and thus be harder to discover by an evolving population. In *Avida* for example, the inclusion of a simple NOT task in the genome requires only one logical NAND instruction (Ofria and Wilke, 2004). On the other hand, the algorithmically complex EQU task requires at least 5 NAND instructions (and many others for proper bookkeeping). Given just this simple restriction, the number of genotypes that encode the NOT task is on the order of a million times larger than genotypes that encode the EQU task. Therefore in this example, the algorithmic complexity of a phenotype is captured in the size of the subset of genotype space in which it is encoded.

From a more information theoretic viewpoint, observation of the phenotype of an organism specifies the set of genotypes that it can possibly occupy. Using these arguments, Wagner (2017) shows that the complexity of an organism's phenotype can be defined as the excess information obtained upon observing a phenotype, which increases as the size of the genetic space decreases (Wagner, 2017). The information content (or complexity) of a phenotype P that contains $|G_P|$ genotypes in it, out of a total of $|G|$ genotypes in the entire space, is then given by,

$$C_P = \log_2 |G| - \log_2 |G_P| \quad (4)$$

This metric is analogous to the definition of biological complexity used previously (Adami et al., 2000). As demonstrated above, it can also be shown to be positively correlated with more functional measures like the algorithmic complexity of tasks that determine the phenotype (Fortuna et al., 2017). The population complexity is calculated as the mean of the phenotypic complexity of all organisms present in the population. The *basal complexity* (C_0) of a host population is the mean population complexity evolved in the absence of parasites. Here, we measure the increase

in complexity upon addition of parasites as the fractional increase over this quantity.

Partitioned Network Structure

Multiple different metrics can be used to quantify the structure of a network. In this article, we use the two most prominent measures that we expect to affect the way the population evolves through the genotype-phenotype map. The first quantity is dependent entirely on the set-like structure of the partition and measures the skew of the genotype-phenotype map. Given a phenotype P_i that has $|G_{P_i}|$ genotypes under it, the skew of the network is given by,

$$\text{skew} = \sum_i \frac{|G_{P_i}|}{|G|} \log_2 \left(\frac{|G_{P_i}|}{|G|} \right) \quad (5)$$

The highest skew (zero) is achieved for partitions that have all the genotypes under a single phenotype. The lowest value ($-\log N_P$) is achieved when the genotypes are distributed equally between the different phenotypes, where N_P is the total number of phenotypes.

The network structure can also be summarized from the perspective of each genotype (node) in the network. Here we look at the heterogeneity in connections (phenotypes that a node connects to) for each genotype node and average it over the entire network to get a measure of neighborhood heterogeneity. Given a node, its connective *neighborhood heterogeneity* (NH) is defined as,

$$\text{NH of a given node} = - \sum_k \frac{E(P_k)}{L} \log_2 \left(\frac{E(P_k)}{L} \right) \quad (6)$$

Where $E(P_k)$ is the number of edges from that node that end up in phenotype P_k and L is the genotype length (the total number of one-step mutants of a L -length bitstring is equal to L). For a given node, this quantity is maximized when it has one-step mutants distributed evenly across all phenotypes and minimized when all its one-step mutants belong to the same phenotype.

Population dynamics on the GP-map

We use the Gillespie method to stochastically model interactions between host and parasite populations on the coarse-grained phenotype map (Gillespie, 1977). Parasites and hosts interact following a matching-alleles model, where only identical phenotypes are allowed to interact. The interaction equations used are as follow,

$$H_i \xrightarrow{r} 2H_i \quad (\text{Host-birth}) \quad (7)$$

$$H_i + P_i \xrightarrow{b} 2P_i \quad (\text{Parasite reproduction}) \quad (8)$$

$$P_i \xrightarrow{c} \emptyset \quad (\text{Parasite death}) \quad (9)$$

$$H_i \xrightarrow{P_{ji}^m} H_j \quad (\text{Host phenotype transition}) \quad (10)$$

$$P_i \xrightarrow{P_{ji}^m} P_j \quad (\text{Parasite phenotype transition}) \quad (11)$$

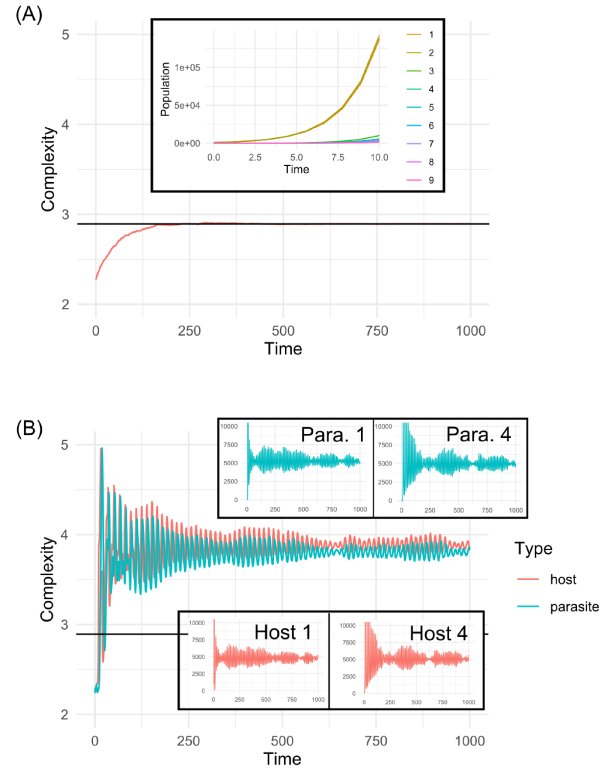


Figure 2: Increase in steady state complexity upon the addition of parasites. For this experiment we used a randomly generated GP-map with a genome length of 10 and 10 different phenotypes. The black line is the basal complexity, C_0 achieved when the host population occupies each genotype with equal likelihood. (A) Host population complexity without parasites relaxes to the basal complexity value over time (Inset: Population dynamics of individual host phenotypes without parasites). (B) Host population complexity in the presence of parasites relaxes to a value greater than the basal complexity (Insets: Population dynamics of host and parasite phenotypes 1 and 4).

Where P_x and H_y are parasites and hosts belonging to x^{th} and y^{th} phenotypes, respectively. P_{ji}^m is the transition rate as obtained from the phenotype transition matrix P at a given mutation rate m . The quantities r , b and c are host-birth, parasite-infection and parasite-death rates respectively. We use the following parameter values in this paper:

Parameter Name	Symbol Name	Value
Host birth rate	r	5.0×10^{-1}
Parasite reproduction/virulence	b	1.0×10^{-4}
Parasite death rate	c	5.0×10^{-1}
Mutation rate	m	2.0×10^{-3}

These values were chosen by simulating the system under a wide range of values and isolating a set of parameters that gave non-zero, finite steady state populations for all the gen-

erated genotype-phenotype maps. Note that due to a lack of either a resource or a space limit for the hosts, the host population increases indefinitely in the absence of parasites (See figure 2A inset). The populations are initialised with 1000 hosts in each of the three lowest complexity phenotypes and 50 parasites of the second lowest complexity phenotype.

Generation of random GP-maps

The GP-maps used in this work were generated using a two-step process (given genotype length L and N_P number of phenotypes),

1. Randomize the ordered set of genotypes (Size 2^L) and randomly place $N_P - 1$ separators between the elements to generate N_P partitions containing different numbers of genotypes.
2. Arrange the partitions in the descending order of number of genotypes to get the set of phenotypes ordered by complexity.

For all the simulations except those in figure 5, we used $L=10$ and $N_P=10$. For figure 5, we started with a specific GP-map ($L=10$, $N_P=7$) that was constructed by hand to have a low node-heterogeneity. This map was then mutated by genotype exchange to give maps with different network structure but the same skew.

Results

Parasites drive escalations in host phenotypic complexity

As seen in earlier work by Zaman et al. (2014), the addition of an initial population of parasites leads to an increase in the steady state complexity reached by the host population (Figure 2b). In the absence of parasites the host population equilibrates over the entire genotype map leading to a basal level of complexity, C_0 (Figure 2a).

Note that both the mean host complexity and phenotypic heterogeneity of the host population changes over time in these simulations (Figure 2b). To get a single steady-state value of these quantities we take an average over the last 100 time-steps of the simulation.

Basal complexity on a genotype-phenotype map is equal to the magnitude of the skew

We find that in the absence of parasites, the host population equilibrates to a distribution such that the phenotypic complexity is exactly equal to the negative skew of the underlying genotype-phenotype map (Figure 3). This occurs because the distribution of host organisms across different phenotypes is expected to be in proportion to the number of genotypes under these phenotypes on a given GP-map. The expression for the skew gives the basal complexity if we assume the distribution of hosts mirrors that of the genotypes (See equation 5).

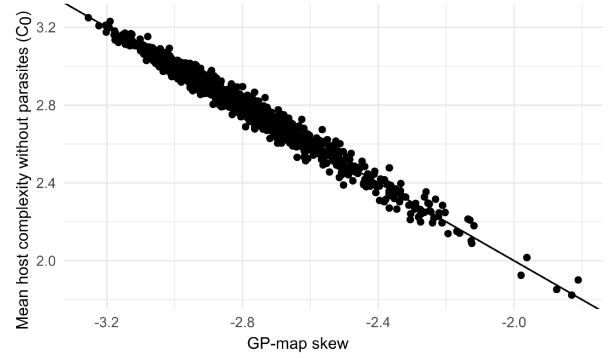


Figure 3: Mean host complexity evolved in the absence of parasites (C_0) as a function of GP-map skew. Note that these two are expected to be exactly equal in magnitude if the host population does distribute proportionally over the entire GP-map. The solid line denotes $y = -x$.

Skew of the GP-map positively affects the evolved phenotypic complexity

By definition, highly skewed GP-maps have a large number of genotypes occupying low complexity (common) phenotypes and few genotypes under high complexity (rare) phenotypes. If parasites are unleashed in a host population equilibrated on such maps, parasites will evolve to reflect a distribution where most of the parasite phenotypes lie under the low complexity phenotypes (i.e., they will target the most common host). In that case, there would be a significant push for host populations to move towards the higher complexity phenotypes because they are scarcely targeted by parasites. At the same time, maps with less skew have basal complexities closer to the level of the most complex phenotypes. Therefore, we hypothesize that an increasingly skewed genotype-phenotype map will lead to greater levels of coevolved complexity when parasites are introduced. We found this to hold true when tested with randomly generated genotype-phenotype maps (Figure 4). Note that we measured the ratio of complexity in the presence and in the absence of parasites, such that what we are measuring is the effect of coevolution on the escalation in complexity over the basal level.

As noted before, the basal complexity over which we calculate this escalation is equal in magnitude to the skew of the underlying map (but opposite in sign). To show that this correlation is not an effect of the lowering of basal complexity for highly skewed maps, we also plot the mean host complexity against the skew of the genotype-phenotype maps (Figure 4 inset). It is interesting to note that for extremely skewed GP-maps, it is possible to get an almost 2.5-fold increase in complexity.

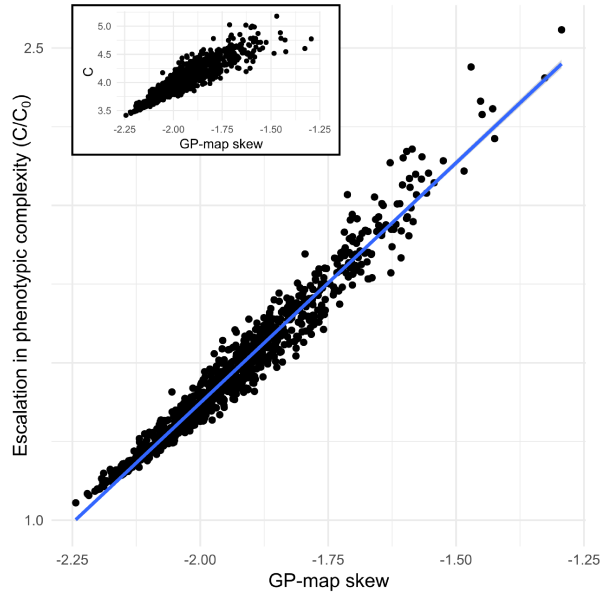


Figure 4: Escalation in host population complexity over basal (C/C_0) plotted against the skew of randomly generated GP-maps. The blue line indicates a linear regression with shaded regions indicating the 95% confidence interval (indistinguishable). Inset: Mean host population complexity (C) plotted against the skew of the genotype-phenotype map. Note that the basal complexity evolved without parasites is the negative of the skew of the genotype-phenotype map.

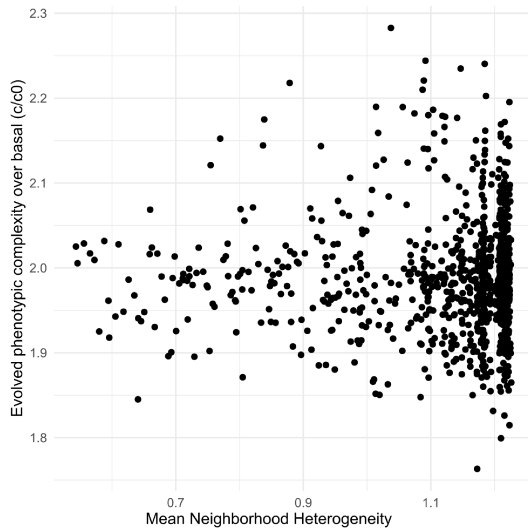


Figure 5: Evolved host population complexity over basal (C/C_0) plotted against the mean neighborhood heterogeneity for randomly generated GP-maps with a constant skew (see text for details on map generation).

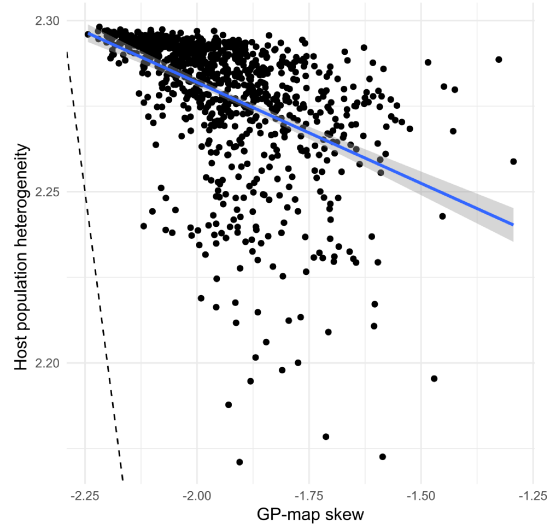


Figure 6: Evolved host population heterogeneity plotted against the skew of randomly generated GP-maps. The dotted black line is the expected population heterogeneity in the absence of parasites (equal to the skew in magnitude). The blue line indicates a linear regression with shaded regions indicating the 95% confidence interval.

GP-maps with higher mean neighborhood heterogeneity have more variable complexity

Neighborhood heterogeneity (defined above) is a property of genotypes that reflects the interconnected nature of the phenotypic space. Certain connections between genotypes may take a sequence from a very low complexity phenotype to the highest complexity phenotype in a single mutation. At the same time, it is possible to have GP-maps where mutations only explore phenotypes with similar levels of complexity. In this respect, this quantity is analogous to the *evolvability* of the GP-map, as increasing neighborhood heterogeneity would allow easier transitions between phenotypes.

To control for skew effects, we generated random maps using a Markov chain starting from a partition with a low mean neighborhood heterogeneity. In each step, we mutate this GP-map by swapping the phenotype labels between pairs of genotypes, thus generating maps with identical skew but different structure. Following this, we performed host-parasite coevolution on these maps to get the steady-state host complexities over the basal level. We find that coevolution using GP-maps with higher mean neighborhood heterogeneity lead to more variable outcomes in complexity (Figure 5). A large neighborhood heterogeneity indicates the presence of extremely evolvable genotypes that might act bidirectionally in a neutral environment – either taking a low complexity genotype to a very high complexity phenotype in the next generation, or vice versa.

Host population diversity at steady state is dependent on the structure of the GP-map

Another interesting – and slightly more tangible – feature of an evolved population is the heterogeneity of the final phenotypes (i.e., diversity) of hosts. In the absence of parasites, the diversity (defined as the Shannon entropy of the phenotype distribution) of the host population mimics the magnitude of the skew of the partition – as the population equilibrates between the phenotypes in proportion to their genotype counts. In the presence of parasites however, we see that the difference between the basal and realized phenotypic heterogeneity increases as the skew of the map increases – indicating that the steady state distribution of phenotypes in the coevolved host population is not just a simple additive effect of parasites increasing diversity above the basal level (Figure 6).

Conclusion

Using a simple model of host-parasite coevolution we show that the escalation in host complexity upon addition of parasites is positively correlated with the skew of the underlying genotype-phenotype map (Figure 4). Although host diversity (i.e., heterogeneity) decreases with GP-map skew, the difference between host diversity in the absence of parasites (Figure 6 – dashed line) and the evolved level of diversity when parasites are included (Figure 6 - blue line) increases. Together, these results suggest that increases in complexity and diversity are driven more strongly by the way antagonistic coevolution is deforming the fitness landscape (the map between phenotypes and reproductive success) than the underlying genetic topology.

While we have presented two major metrics that affect the outcome, their relationship with the final complexity is in no way explicit. In addition, our simulations are performed at specific mutation levels, virulence, and birth/death rates; it is likely that the relationship between these variables and the final complexity also changes under different GP-map structures – an area of enquiry that has direct implications for work that aims to characterize host-parasite population dynamics more generally.

Specifically, our coevolution experiments occurred in an otherwise neutral fitness landscape. How robust these results are to more complex genotype-phenotype maps is one interesting direction to explore. For example, does having genotypes that represent inviable organisms change the relationship between neighborhood heterogeneity and complexity? We also assume that the hosts and parasites have similar genotype-phenotype maps, an assumption not generally true for most systems observed in nature. In the future, we seek to explore the effect of variation in the genotype-phenotypes maps of both hosts and parasite independent of each other.

At the values of the parameters used in these simulations, we do not see a strong dependence on the local network structure of the GP-map, other than the fact that networks

with higher neighborhood heterogeneity seem to be more variable (Figure 5). Perhaps our requirement of coexistence across the entire range of random GP-maps was too stringent. It might have constrained parameter values to regions where the dynamics depend primarily on skew (a statistical property of the distribution over genotypes) rather than being dependent on the network structure of the map between genotypes and phenotypes.

The broad goal of this article is twofold – firstly, to highlight that we can improve the predictive capabilities of artificial life frameworks by studying the structure of the underlying mapping between genotypes and phenotypes, and secondly, to demonstrate that coevolutionary dynamics on differing genotype-phenotype mappings is an interesting area for theoretical study. The structure of such systems is similar to the two-level classification of indistinguishable micro-(genotypes) and distinguishable macrostates (phenotypes) in classical statistical physics, except the system is not in equilibrium due to the presence of external sources and sinks of energy (manifesting here as the differential birth and death processes). While expected, it is interesting to see that the outcome of a dynamical process on such a map is so strongly tied to its statistical properties. It might thus be possible to draw accurate predictions about coevolutionary populations based entirely on the genotype-phenotype map they evolve on – especially in regimes where neutral and selective forces are comparable in magnitude.

The major implication of these results is that certain features seen in coevolved host populations – like escalation in diversity, complexity, and the effects of other parameters on these – might be outcomes sensitive to the statistical features of the genotype-phenotype map. These properties, being affected by the architecture of the underlying mapping between the genotype and phenotype space, thus conform to the idea of evolutionary spandrels as discussed extensively in classical and recent works in evolutionary theory (Gould et al., 1979; Valverde et al., 2018). The extent to which biological processes themselves shape these structures remains an open and exciting question. For example, on multidimensional fitness landscapes, parasites may evolve to occupy regions with higher neighborhood heterogeneity, which could improve the evolutionary access to several host genotypes. In the context of artificial life, we believe that accounting for these structure-dependent outcomes will greatly increase the predictive power of Alife frameworks.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. DEB-1813069. BK would like to acknowledge the Department of Science & Technology, Government of India for providing the KVPY fellowship (SA-1410015).

References

- Adami, C., Ofria, C., and Collier, T. C. (2000). Evolution of biological complexity. *Proceedings of the National Academy of Sciences*, 97(9):4463–4468. Publisher: National Academy of Sciences Section: Physical Sciences.
- Floreano, D. and Nolfi, S. (1997). God Save the Red Queen! Competition in Co-Evolutionary Robotics. Conference Name: 2nd Conference on Genetic Programming Number: CONF.
- Fortuna, M. A., Zaman, L., Ofria, C., and Wagner, A. (2017). The genotype-phenotype map of an evolving digital organism. *PLOS Computational Biology*, 13(2):e1005414. Publisher: Public Library of Science.
- Gillespie, D. T. (1977). Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry*, 81(25):2340–2361.
- Goldsby, H. J., Knoester, D. B., Ofria, C., and Kerr, B. (2014). The Evolutionary Origin of Somatic Cells under the Dirty Work Hypothesis. *PLOS Biology*, 12(5):e1001858. Publisher: Public Library of Science.
- Gould, S. J., Lewontin, R. C., Maynard Smith, J., and Holliday, R. (1979). The spandrels of San Marco and the Panglossian paradigm: a critique of the adaptationist programme. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 205(1161):581–598. Publisher: Royal Society.
- Hillis, W. D. (1990). Co-evolving parasites improve simulated evolution as an optimization procedure. *Physica D: Nonlinear Phenomena*, 42(1):228–234.
- Miikkulainen, R. and Stanley, K. O. (2004). Competitive Coevolution through Evolutionary Complexification. *Journal of Artificial Intelligence Research*, 21:63–100. arXiv: 1107.0037.
- Nelson, C. W. and Sanford, J. C. (2011). The effects of low-impact mutations in digital organisms. *Theoretical Biology & Medical Modelling*, 8:9.
- Ofria, C. and Wilke, C. O. (2004). Avida: a software platform for research in computational evolutionary biology. *Artificial Life*, 10(2):191–229.
- Valverde, S., Piñero, J., Corominas-Murtra, B., Montoya, J., Joppa, L., and Solé, R. (2018). The architecture of mutualistic networks as an evolutionary spandrel. *Nature Ecology & Evolution*, 2(1):94–99. Number: 1 Publisher: Nature Publishing Group.
- Wagner, A. (2017). Information theory, evolutionary innovations and evolvability. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 372(1735):20160416. Publisher: Royal Society.
- Wagner, A. P., Zaman, L., Dworkin, I., and Ofria, C. (2020). Behavioral Strategy Chases Promote the Evolution of Prey Intelligence*. In Banzhaf, W., Cheng, B. H., Deb, K., Holecamp, K. E., Lenski, R. E., Ofria, C., Pennock, R. T., Punch, W. F., and Whittaker, D. J., editors, *Evolution in Action: Past, Present and Future: A Festschrift in Honor of Erik D. Goodman*, Genetic and Evolutionary Computation, pages 225–246. Springer International Publishing, Cham.
- Watson, R. A. and Pollack, J. B. (2001). Coevolutionary dynamics in a minimal substrate. In *Proceedings of the 3rd Annual Conference on Genetic and Evolutionary Computation, GECCO'01*, pages 702–709, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Wilke, C. O., Wang, J. L., Ofria, C., Lenski, R. E., and Adami, C. (2001). Evolution of digital organisms at high mutation rates leads to survival of the flattest. *Nature*, 412(6844):331–333.
- Zaman, L., Meyer, J. R., Devangam, S., Bryson, D. M., Lenski, R. E., and Ofria, C. (2014). Coevolution Drives the Emergence of Complex Traits and Promotes Evolvability. *PLOS Biology*, 12(12):e1002023. Publisher: Public Library of Science.