

RESEARCH ARTICLE

Lexical Landscapes as large *in silico* data for examining advanced properties of fitness landscapes

Victor A. Meszaros[☉], Miles D. Miller-Dickson[☉], C. Brandon Ogbunugafor[✉]*

Department of Ecology and Evolutionary Biology – Brown University, Providence, Rhode Island, United States of America

☉ These authors contributed equally to this work.

* brandon_ogbunu@brown.edu

Abstract

In silico approaches have served a central role in the development of evolutionary theory for generations. This especially applies to the concept of the fitness landscape, one of the most important abstractions in evolutionary genetics, and one which has benefited from the presence of large empirical data sets only in the last decade or so. In this study, we propose a method that allows us to generate enormous data sets that walk the line between *in silico* and empirical: word usage frequencies as catalogued by the Google ngram corpora. These data can be codified or analogized in terms of a multidimensional empirical fitness landscape towards the examination of advanced concepts—adaptive landscape by environment interactions, clonal competition, higher-order epistasis and countless others. We argue that the greater *Lexical Landscapes* approach can serve as a platform that offers an astronomical number of fitness landscapes for exploration (at least) or theoretical formalism (potentially) in evolutionary biology.

OPEN ACCESS

Citation: Meszaros VA, Miller-Dickson MD, Ogbunugafor CB (2019) *Lexical Landscapes* as large *in silico* data for examining advanced properties of fitness landscapes. PLoS ONE 14(8): e0220891. <https://doi.org/10.1371/journal.pone.0220891>

Editor: Josh Bongard, University of Vermont, UNITED STATES

Received: May 24, 2019

Accepted: July 25, 2019

Published: August 12, 2019

Copyright: © 2019 Meszaros et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Code and data sets featured in this manuscript can be found on Github: https://github.com/OgPlexus/Lexical_Landscapes.

Funding: CBO was funded by NSF RII Track-2 FEC (Award Number: 1736253).

Competing interests: The authors have declared that no competing interests exist.

Introduction

Historically, theoretical population genetics has often progressed by using contrived or even artificial systems and models to recapitulate the properties of natural, real-world populations of organisms [1, 2]. The fitness landscape is one of the triumphs of theoretical population genetics, an abstraction that changed how evolutionary biologists studied the process of adaptive evolution [3], and one whose original iterations were entirely theoretical. While these models are always engineered with a set of necessary constraints, they have been central to the growth of theory in modern population genetics [4].

In recent years, fitness landscapes have benefited from advances that enabled the use of empirical data towards the construction of empirical fitness landscapes [5]. Combinatorial data sets—where suites of mutations are engineered in every possible permutation—are the gold standard for these types of studies [6]. They were critical for introducing the concept of the adaptive trajectory, and have since been used as an innovation space for methods to detect higher-order epistasis [7], for metrics to calculate the speed of adaptive evolution [8], and for

more rigorous attempts at predicting or steering evolution [9, 10]. The limitations of combinatorial data sets are that they tend to only focus on suites of mutations within a single gene of interest, and that there are relatively few such data sets in existence [11–14]. Regardless of source, fitness measurements for these landscapes are often taken in a small number of environments, which limits our understanding of how the effect of mutations might be affected by environments. Note that this is even a problem for existing studies featuring simulated and *in silico* fitness landscapes [15, 16].

In this study, we propose an instrument—*Lexical Landscapes*—for generating theoretical fitness landscapes, that is not based on an algorithm. Instead, it is built on an open-source and well-established data set that can be easily analogized as fodder for the study of advanced topics in theoretical population genetics: Google Books ngram data corpora. In a prior study, this concept was introduced as a model for evolution through “protein space,” as it builds on a highly effective analogy authored by John Maynard Smith [17, 18], and can serve as a means of teaching and communicating concepts in evolutionary genetics as well. Here we double-down on this idea by arguing that the utility of the Google Books ngram corpora is not only pedagogical but also exploratory and scientific: one can use this data set to test and generate hypotheses, and develop theory in modern evolutionary genetics that exceeds the reach of current data sets (empirical and simulated). Even more, the accessibility of the data set, and connection to common words makes it easy to codify, discuss, and cross-reference.

We first outline the specific data science and computational methods necessary to generate a set of *Lexical Landscapes* of a certain kind. We then demonstrate the utility of these sets by exploring how *Lexical Landscapes* can recapitulate several standard and advanced properties of evolutionary genetics on combinatorial data sets, such as fitness landscape topography and the accessibility of adaptive trajectories. For example, we introduce an “environmental context” analogy to these data, which allows us to rigorously compute properties of fitness landscapes across environments. We then move onto the elusive concept of higher-order epistasis, and examine how it is affected by environmental context. Summarizing, we re-emphasize the breadth of concepts that can be explored with *Lexical Landscapes* and speculate on its potential as an instrument for modern theoretical biology.

Methods

Using Google ngram values to generate empirical fitness landscapes: Conceptual challenges

A prior study introduced the use of Google ngram data as a pedagogical and communicative tool for evolutionary genetics [18]. In this study, we expand this idea and argue that the Google ngram corpora has utility for thinking about more advanced concepts in evolutionary biology. To effectively utilize *Lexical Landscapes*, it is critical that one potentially confusing idea is fully clarified: while many of the patterns we observe in the data set might be reflective of cultural [19] or evolutionary linguistic phenomenon, *Lexical Landscapes* are not engineered to study the evolution of language. Rather, they offer a transparent, open-access reservoir of data that can be easily translated into a form similar in structure to other biological and *in silico* data sets used to generate fitness landscapes. In addition, the *Lexical Landscapes* approach in this study is confined to small, combinatorial-style fitness landscapes. While one can imagine other uses for the data examined in this study (e.g. larger landscapes), we have rooted our study in existing studies and in the original John Maynard Smith conceptualization [17, 18]. In [Table 1](#), we define the concepts and terms of *Lexical Landscapes* and their evolutionary analogues.

Table 1. Defining concepts and terminology in Lexical Landscapes.

Concept	Definition	Example
<ul style="list-style-type: none"> • Fitness landscape • Adaptive landscape • Fitness Graph 	These terms are used to describe a kind of genotype-phenotype map that organizes information in a manner that communicates details about the evolutionary process.	See Fig 2
<ul style="list-style-type: none"> • Accessible pathway or trajectory 	An evolutionary pathway through a fitness landscape where consecutive step leads to an increase in the fitness proxy.	See Fig 3
<ul style="list-style-type: none"> • Within path competition (C_w) 	Across a given accessible pathway there is competition between alleles. Recent work has demonstrated that the this total competition is powerfully associated with the speed of evolution across a given pathway	See Fig 5
<ul style="list-style-type: none"> • Environment or context 	How evolution occurs is profoundly influenced by its environmental context. In Lexical Landscapes, either the time dimension (continuous variable) or language variable (categorical) can be used as an analogy for environment or context.	Time or language. For example, the CARS → SOME landscape can be constructed at any of the available time points in which there are data. In Fig 2, we construct landscapes
<ul style="list-style-type: none"> • Genotypic Context 	How a fitness landscape for a given gene is affected by the whole-organism genome in which it is embedded. For example, in transgenic studies where a gene from organism A is engineered into the genome of organism B. It is possible that this change in genomic context can change the phenotypic effects that gene. In this study, we use British-English as an analogy. The CARS → SOME landscape features can be calculated in certain sub-“genera” of English, like British-English.	See S1 File. Supplemental Appendix

<https://doi.org/10.1371/journal.pone.0220891.t001>

We first outline a method for collecting and curating these data (Fig 1). We then put these data to use through various calculations and simulations, which highlight the cutting-edge evolutionary questions that can be interrogated with this data set.

Data acquisition and curation

The strength of the *Lexical Landscape* approach resides in its expansiveness. Given the *English* alphabet of 26 letters, there are over 600 different 2-letter “alleles” (different word variants) and over 17,500 3-letter alleles. While the majority of words will have near-zero usage frequencies, or what we call *Lexical Fitness*, there might be utility in studying a large portion of these landscapes. Thus, we propose that the innovation of this approach resides in its ability to create an enormous number of non-arbitrary fitness landscapes. In order to demonstrate how this

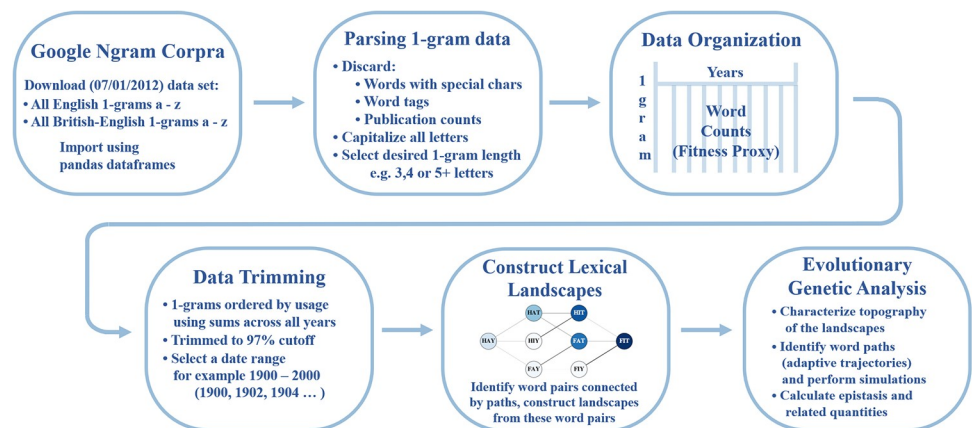


Fig 1. Methods flowchart. The flowchart above describes the methodology used to construct *Lexical Landscapes* and carry out several of the additional analyses discussed in this study.

<https://doi.org/10.1371/journal.pone.0220891.g001>

data can be leveraged for evolutionary theory, we have focused on combinatorial and empirical landscapes of smaller size. A fuller description of the size of the *Lexical Landscape* data set (actual and curated) can be found in the Supplementary Appendix.

To construct the data set of *Lexical Fitnesses*, we formatted data provided in Google's ngram corpora. Specifically, we downloaded the entire *English* corpora of *1-grams*—or single words, as opposed to 2-grams for instance, which are pairs of words—for all the letters in the alphabet. We discarded 1-gram data sets containing numbers or special characters. The data from Google's corpora originally included “word tags” (e.g. __NOUN__, __VERB__, __ADJ__) specifying the grammatical context in which the word was used—this provided a breakdown of how often a given word was used as a noun, adjective, etc. In our data set, we simply removed these tags and ignored their grammatical context. The word counts associated with each word are therefore the *total* usage counts, which are sums over these grammatical contexts. Lastly, the original data includes counts enumerating the number of books each word appeared in; these were also discarded for our purposes (although, this information and others that are available could be useful for other purposes).

For this study, the data is composed of all 3, 4, and 5-letter words in Google's *English* corpora, along with the number of times each word was used in every other year beginning in 1900 and ending in 2000. The ngram word frequencies are taken from books that Google has been able to survey (approximately 6% of all books ever published [20]). We limited the data to every other year to expedite computations and we chose to use the 20th century data as it represents the most modern full century from the Google corpora—it is also the most densely populated in terms of word usage. Note, however, that the Google ngram corpora data go back to the 1500s, and so there is nothing preventing the construction of *Lexical Landscapes* for any of these years.

We divided the data into the three word-lengths—3, 4, and 5—and within each we ordered the words according to word popularity. More precisely, for each word within each word-length category, we summed the total word count usage for that word across all the years in our data set—a total of 51 years: 1900, 1902, . . ., 1998, 2000. Words with a higher sum were assigned a higher popularity. For instance, we found that the most popular 3, 4 and 5-letter words—calculated in this manner—were ‘THE’, ‘THAT’, and ‘WHICH’, respectively. In order to use words that are fairly popular and to avoid acronyms and initialisms—also to facilitate quicker calculations—we truncated each list of words down to the top 200 most popular 3-letter words, the top 1500 most popular 4-letter words, and the top 5,000 most popular 5-letter words. The popularity cutoffs were chosen so that the sum of all word counts in each truncated list represented roughly 97% of the total summed usage-count of all words in that word-length category. We chose 97% so that the data set was as large as possible without including many of the very numerous, though infrequently used, acronyms/initialisms.

Within each truncated list, we identified various *word-paths*. A word-path is characterized by a starting and ending word, and comprises a collection of words connected by way of changing one letter at a time. More precisely, a letter-change occurring at some location in the word is a swap of a letter in the beginning word for the corresponding letter in the ending word, such that at each step the combination of letters is a word itself. An example can be seen in Fig 2. We will call the collection of *all* letter combinations between two terminal words a *landscape*—we will analogize this to a *fitness* landscape—and two examples of which are shown in Fig 3. Word-paths were interesting to consider since the words that form a path each have some kind of lexical significance, and can be analogized with *genotypes* with sufficiently high fitness. Whereas, arbitrary combinations of letters—which do not make good *lexical* sense—are analogized with genotypes that do not fair well in biological settings, having low fitness (i.e. they do not make good *biological* sense). Using word-usage frequency as a proxy for

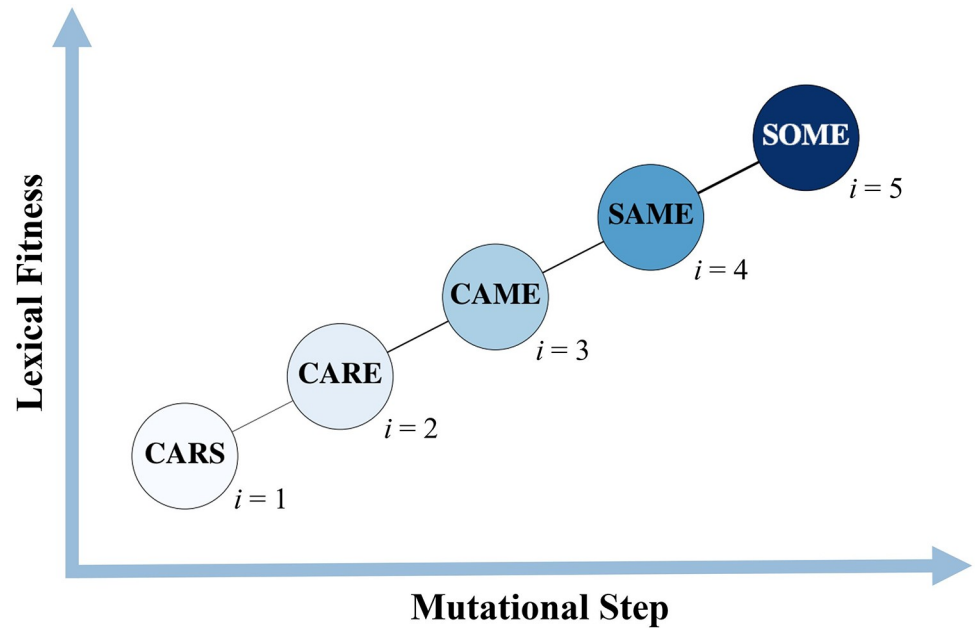


Fig 2. Adaptive trajectory. A hypothetical “uphill” word-trajectory in the CARS → SOME landscape. Each word is indexed by i , referenced in Eq 1, and the color gradient indicates the fitness value r_i , with darker blues showing higher fitness. The edges connecting the nodes are thicker and thinner, depending on the difference in the fitness values (in absolute value).

<https://doi.org/10.1371/journal.pone.0220891.g002>

fitness, we identify among the set of word-paths those which are “evolutionarily favorable”—i.e. those for which the fitness value increases along the path, in a given year. We refer to these paths as *accessible* or sometimes *uphill* paths. The details of the algorithm used to identify these paths, and how to access the Python script, are described in the Supplementary Appendix. It is very important to note that even with the selective criteria used to curate the greater ngram corpora, we can still generate over 1 million total fitness landscapes, that is, combinatorial sets with *Lexical Fitness* values between pairs of alleles. This would constitute, by many measures, the largest set of fitness landscapes in existence.

Choice of example Lexical Landscapes used for illustrative purposes

In order to examine various properties of fitness landscapes using *Lexical Landscapes*, we chose two model landscapes in the set of 4-letter words, CARS → SOME and POEM → LAST.

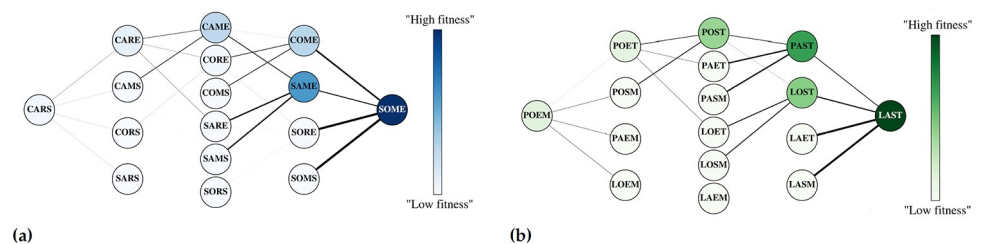


Fig 3. Fitness graphs for the CARS → SOME and POEM → LAST landscapes. Nodes are connected if they differ by a single letter. Darker nodes represent higher fitness values, and edges are weighted by the difference in the fitness values (in absolute value).

<https://doi.org/10.1371/journal.pone.0220891.g003>

These two were chosen among thousands of possible landscapes because they carried features that make them useful examples:

- Their topography changed with environmental contexts, which made it likely that evolutionary dynamics might differ across context.
- Their topography is rugged, indicating the presence of epistasis. Rugged landscapes often create non-intuitive evolutionary patterns.
- They contained multiple paths that were accessible across several environments. These are also features of landscapes that are likely to offer non-intuitive dynamics, as these landscapes possess multiple accessible paths which can *compete* in an evolutionary sense.

Note that none of these above outlined features of model landscapes are essential for the use of *Lexical Landscapes* as an instrument for studying and exploring properties of fitness landscapes. We only imposed these criteria on the landscapes in order to illustrate the potential utility of this tool. Any one of thousands of *Lexical Landscapes* might be generated for the study of advanced properties of fitness landscapes. In the Supplementary Appendix, we provide similar analyses for 3 and 5-letter word *Lexical Landscapes*.

In what follows, we show how three concepts in evolutionary biology can be explored with this approach. We first consider how fitness values can be used as a proxy for growth rates in order to simulate the “growth” and evolution of the various alleles, or word combinations, that comprise a landscape. We compare simulated dynamics in various contexts. Second, we explore how lexical fitness values can be used to construct a metric for the expected *time* associated with evolution along a particular trajectory or *word-path*—we refer to this metric as the within-path competition. Lastly, we use lexical fitness to calculate the degree of epistasis present in the data set. We observe how this degree may vary over years across the 20th century.

Evolutionary dynamics with the Discrete Asexual Reproductive Population Simulator (DARPS)

Here we discuss the simulated evolutionary dynamics for the two example 4-letter landscapes (CARS → SOME and POEM → LAST). The purpose is to illustrate how simulations across *Lexical Landscapes* can demonstrate many properties, simple and advanced, that are direct analogues to biological processes.

Each landscape consists of the sixteen letter combinations for a given word pair, featuring all possible letter-swaps between the words in the word pair (Fig 3). In the simulation, the *population* of the first word (allele) in the word-pair was set to 1000, while all other words had an initial population of 0. Each population was allowed to evolve at each time step according to some fixed probability of mutation (we chose 10^{-8}), a mutation rate that is on the order of what we observe in microbial populations. Each word was also assigned a growth rate which was correlated with its *Lexical Fitness* for a given year in the following way: we first calculated the average fitness value of all words in the landscape for the given year. Then, we use the fitness values in the landscape divided by this average for the given year as the growth rates for each word. Note that the growth rates were fixed throughout the simulation and only depended on the fitness values for the year chosen. In this way, the growth rates are defined by their *relative* fitness (relative to the mean fitness of all words in a given landscape for a given year). We used the *Discrete Asexually Reproducing Population Simulator* (DARPS), a simulator of evolution in large populations of organisms that resembles microbial populations [8, 21]. At each time step in the simulation, or generation, a certain proportion of each word’s population undergoes replication (in the generic exponential sense, with the growth rate as the Malthusian

parameter). Mutations occur during replication according to the probability of mutation. As the simulation progresses, different alleles can rise and fall in frequency. We can visualize the dynamics of these simulations by graphing the fraction of each allele in the population (Fig 4).

Within-path competition and the speed of adaptive evolution across a fitness landscape

We now consider how *Lexical Landscapes* can be used to examine advanced concepts that have never before been explored at the scale that *Lexical Landscapes* offers. A study [8] introduced a term that correlates powerfully with the time associated with evolution across a landscape and is calculated for a given *path*. It is termed the *within-path* competition (C_W) and it is defined by the formula,

$$C_W = \sum_{i=1}^{N-1} \frac{1}{r_{i+1} - r_i} \tag{1}$$

where r_i is the *growth rate* of the *i*th allele, and where N is the length of the evolutionary-path (the number alleles in the path)— C_W can therefore be thought of as a sum over the *links* (or edges in a graph, such as in Fig 3) in a path. C_W is typically computed using growth rate measurements, however, in this work we use the *Lexical Fitness* values as a proxy for growth rate. When C_W is high, evolution across a trajectory is expected to be slow, and evolution is fast when C_W is low, provided that $r_{i+1} > r_i$ for each i , which assumes that the path is an *uphill* path. In the context of *Lexical Landscapes*, we calculate C_W along *word-paths* in the landscape. These paths are analogous to a series of mutational steps in a 4 loci allele, progressing along an evolutionary path, such as 0000 → 0001 → 0011 → 0111 → 1111. Fig 5 shows two such word-paths, one in each of our example landscapes. One can observe how drastically this value can vary across time. We present more examples of C_W in the Supplemental Appendix. For a more rigorous technical treatment of the topic, we point readers to the reference where the C_W was introduced: Ogbunugafor and Eppstein, 2017 [8].

Calculating higher-order epistasis

Epistasis remains a cutting-edge topic in evolutionary biology that continues to be the object of study for a variety of reasons, and measured using diverse methods [13, 22–24]. For our purposes, we use a Walsh-Hadamard transformation of the fitness values, scaled by an additional diagonal matrix, as presented in Poelwijk et al. [22]. We summarize the approach below.

For a given year, the *Lexical Fitness* values for each letter combination in a given landscape (CARS → SOME or POEM → LAST) are arranged into a vector x —for 4-letter words there are 16 fitness values, one for each letter combination in the landscape. In short, this vector will be multiplied by a 16 x 16 square matrix; we then take the absolute value of the output and normalize. The 16 x 16 matrix is the product of a diagonal matrix V and a Hadamard matrix H .

These matrices are defined recursively by,

$$V_{n+1} = \begin{pmatrix} \frac{1}{2}V_n & 0 \\ 0 & -V_n \end{pmatrix}, \quad V_0 = 1$$

$$H_{n+1} = \begin{pmatrix} H_n & H_n \\ H_n & -H_n \end{pmatrix}, \quad H_0 = 1$$

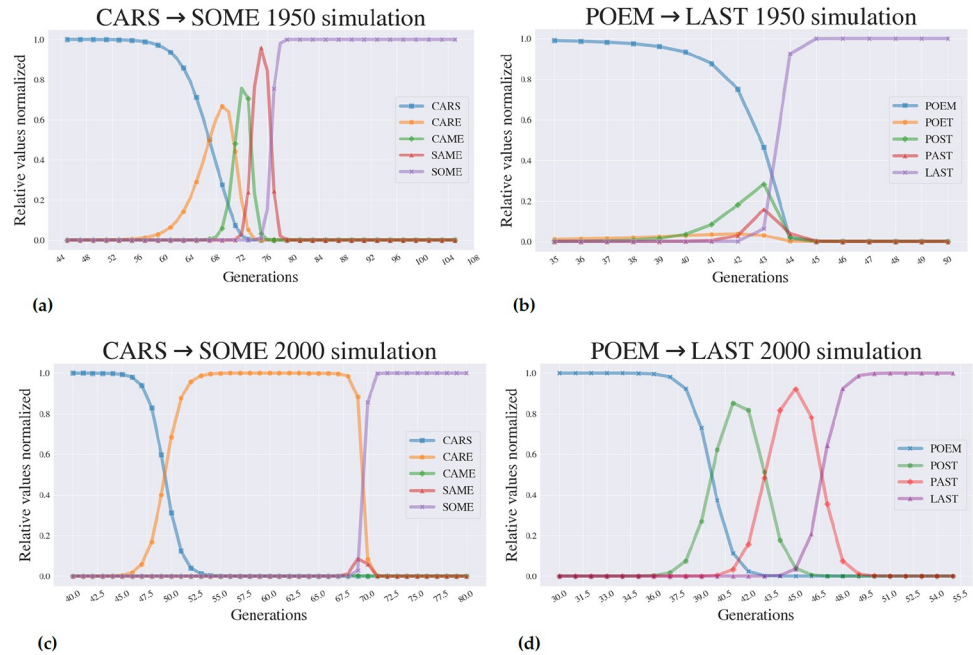


Fig 4. Simulated dynamics across different contexts demonstrate different dynamic properties. Fig 4 (a) and (c) show illustrative simulations across the CARS → SOME *Lexical Landscape*. Fig 4 (b) and (d) show two for the POEM → LAST simulation dynamics. The y-axis shows the percentage of the population occupied by a given word. All simulations begin with 100% of the population fixed at the 0000 (CARS or POEM) genotype and ultimately reach fixation at the 1111 (SOME or LAST) genotype. Note, however, that the dynamics through which this occurs changes as a function of context, or year in our case.

<https://doi.org/10.1371/journal.pone.0220891.g004>

where n is the number of loci ($n = 4$ for our purposes as we consider 4-letter words). The output y of this matrix multiplication is given explicitly by,

$$y = VHx \tag{2}$$

where V and H are the matrices above for $n = 4$. We then take the absolute value of the entries in y and divide each entry by the sum of the absolute values to normalize. In this paper, we

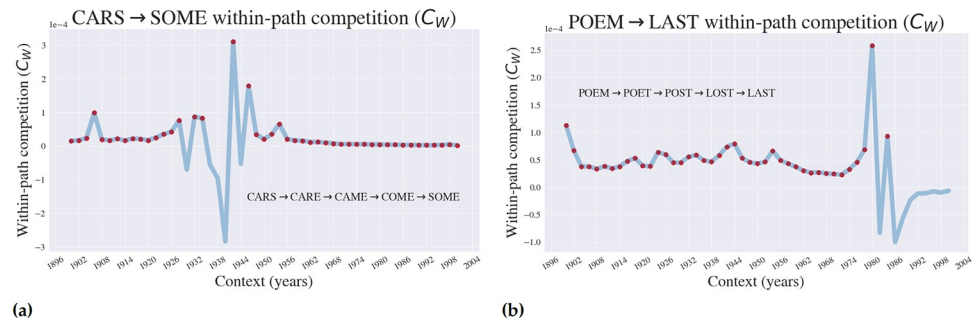


Fig 5. Within-path competition (C_W), a proxy for the speed of evolution changes with context. The *within-path* competition coefficient is shown for specific paths (shown in the figures) in the CARS → SOME (a) and POEM → LAST (b) landscapes. Red dots indicate the years for which C_W was calculated on an *uphill* path, which is to say that the associated path had increasing fitness values as the path is traversed.

<https://doi.org/10.1371/journal.pone.0220891.g005>

refer to these as the *epistatic* coefficients E_i .

$$E_i = \frac{|y_i|}{\sum_j |y_j|} \quad (3)$$

The absolute value and normalization was performed in order to focus exclusively on the magnitude of the epistatic effects. We will sometimes use bit strings such as ‘0101’ as the index i when referring to the epistatic coefficients E_i . For instance, E_{0101} represents a measure of the epistatic effect of “mutations” in the 2nd and 4th letter, or locus.

In addition, one may consider averaging these epistatic coefficients E_i within what we call an *order*, enabling comparisons between orders. An *order* is the collection of mutations of an allele with the same *number* of mutations, or bit-flips (or letter *swaps* in a given combination of letters in the landscape). For instance, the coefficient E_{0000} belongs to what we refer to as the “0th” order (zero letter swaps), whereas E_{0001} , E_{0010} , E_{0100} , and E_{1000} all belong to the “1st” order (one letter swap), etc. Note that, like the 0th order, there is only one coefficient contained in the 4th order, namely E_{1111} . Taking the average of the epistatic coefficients within each order presents information about the general epistatic effect based solely on the number (or order) of mutations (letter swaps) in a given landscape, one swap, two swaps etc. In Fig 6, we present these average epistatic effects for our two landscapes (we label them with the term “absolute mean” since we incorporated absolute values and averages in the calculation). In the Supplemental Appendix, we show all the epistatic orders (without averaging) for our two case landscapes, as well as two additional ones within the 3 and 5-letter categories—in the Supplemental Appendix, we refer to these plots as the *dis-aggregated* epistasis, in contrast to the *aggregated* epistasis we present here in the main text. As past studies have focused on comparing higher-order effects [7], one can glean a lot of information from aggregating all effects by their order. For example, we can observe whether a given landscape is dominated by epistatic effects of a certain order, and speculate as to why this is so.

Results and discussion

Simulations of evolution across *Lexical Landscapes*: Standard and non-standard dynamics

Fig 4 demonstrates that *Lexical Landscapes* can be used for the construction of simulated adaptive evolution as observed in several studies of evolution across empirical fitness landscapes [5, 6, 25, 26]. The *Lexical Landscapes* used as models in this study also demonstrate advanced properties of evolution. While Fig 4a displays the standard step-wise evolutionary trajectory whereby the population reaches appreciable frequencies at all alleles in a given path, Fig 4b–4d demonstrate how certain simulations display features of *stochastic tunneling*, where evolution appears to “skip steps”. That is, when an intermediate genotype makes no appreciable appearance in population space. Prior studies have revealed that stochastic tunneling can happen when populations sizes are large and mutations rates are high. [27, 28].

The accessibility of pathways and within-path competition (C_w ; a proxy for the speed of evolution) changes as a function of context

Fig 5 shows how one may consider visualizing the *time* that evolution takes within a landscape across varying contexts. This translates to the idea that the speed of adaptive evolution is a function of the environment that it’s in. In the Supplemental Appendix, we compare the plots in Fig 5 to two other paths within each landscape, in addition to considering C_w for other word lengths.

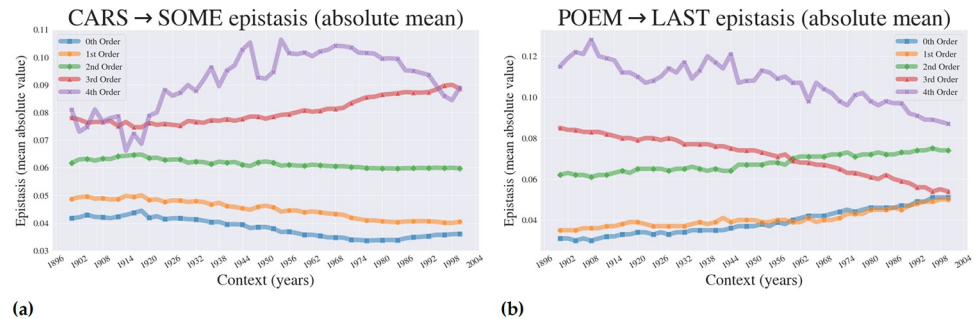


Fig 6. Higher-order epistasis across context. Fig 6 (a) and (b) show how epistatic effects for our example word landscapes can vary across environment, or in this case across time.

<https://doi.org/10.1371/journal.pone.0220891.g006>

The magnitude of higher-order epistasis changes across environmental contexts (years)

Next, we turned our attention to a popular concept in evolutionary genetics: higher-order epistasis. Using the Walsh-Hadamard transformation (see: [Methods](#)), we calculated the aggregated higher-order epistasis for each order (0th–4th) and plotted it across the contexts from 1900–2000. Strikingly, we observe powerful context dependence of higher-order epistatic effects ([Fig 6](#)), for both examined landscapes (CARS → SOME and POEM → LAST). Interestingly, we can also see that across most of the queried environmental contexts (1900–2000), fourth order effects are predominant, indicating that within these landscapes, words derive their utility (in terms of word-frequency, or *Lexical Fitness*) from interactions between all four of the letters. One also notices that the predominant epistatic influence across the landscape changes across contexts. There are, for example, small contextual windows (years) in the CARS → SOME landscape where the 3rd order effects dominate. Keep in mind that the analysis in [Fig 6](#) obfuscates the direction of epistasis: as outlined in the methods, absolute values of the output vectors were computed, which means that information on the sign of individual effects is lost. The data for the entire calculation, however, can be found on GitHub: github.com/ogplexus/LexicalLandscapes.

Conclusion

In this study, we introduce a method for creating a database for a large number of empirical fitness landscapes based on the ngram corpora. Though these data are not “biological,” they are empirical in the sense that the data for the individual nodes arise from a measurement of a natural system such as a language, and are not simulated. For this reason, *Lexical Landscapes* provide an ideal hybrid empirical-*in silico* data set for readily generated fitness landscapes at a size and scope that far exceed that of current technology-limited biological systems.

We also demonstrate how to put these data to use in a way that highlights how one can observe advanced and higher-order phenomena in evolutionary genetics. Importantly, we emphasize that a strength of the *Lexical Landscape* is in its capacity to provide analogies for varying environmental contexts, using time as one such means for variation. We do so using three different exercises.

Firstly, we demonstrate how one can observe evolutionary simulations across these landscapes. These simulations follow the step-wise process of adaptive evolution that has been observed for similarly constructed empirical fitness landscapes. We also observe, however, circumstances where evolution displays non-standard properties. Specifically, the topography of some landscapes displays features of “stochastic tunneling,” where a population appears to

“skip steps” in adaptive evolution. The reality is that evolution isn’t “skipping” steps at all, but rather, that an intermediate allele is present in low frequency, but due to high overall population size and mutation rate, subsequent steps are traversed.

Next, we examined how within-path competition (C_w), a measure of the clonal competition acting along a particular adaptive trajectory. This has been demonstrated to be strongly correlated with the speed of evolution across a certain trajectory [8], and so this analysis offers a demonstration of how one can gain a picture for the predicted speed of adaptation as a function of a fitness landscape occupying some specific niche.

Lastly, and most provocatively, we demonstrate how the magnitude of higher-order epistasis is altered by context. Despite the fact that epistasis is a powerfully controversial topic, very few studies have explored how context influences higher-order epistasis [29–31]. Using *Lexical Landscapes*, we are easily able to generate fitness landscapes, calculate epistatic coefficients within single environments, and demonstrate how those coefficients manifest across a wide number of environments. This is a understudied phenomenon, one for which no general theory or rules have yet been constructed. Using *Lexical Landscapes*, we now have insight that can be applied to existing biological data. There are many reasons why more careful interrogation of how epistatic effects change with environment might be relevant, as it could explain difficulties in recovering individual SNPs of large effect in large-scale genomic studies.

The empirical-*in silico* dichotomy is one that has been a part of population and evolutionary genetics since its inception, and will likely always have a place in evolutionary theory. Certain questions will always benefit from volumes of data that are beyond the scope of what the biological world offers. With *Lexical Landscapes*, we offer a method for generating data sets that can be used to explore many features of adaptive evolution.

In closing, we want to emphasize that the outlined approach is hardly the only way that ngram data might be utilized towards asking questions related to the fitness landscape. We encourage others to expound upon these methods and data, towards studying several yet unseen phenomenon in evolutionary theory.

Supporting information

S1 File. Supplemental appendix. This contains a number of important explanations and extensions of the *Lexical Landscape* approach. They include:

- A description of the size of the data set utilized in this manuscript
- An examination of how *Lexical Landscapes* can be explored in other language subsets (e.g. British English)
- Further details and calculations of higher-order epistasis
- Added investigations into landscapes of 3 and 5 letter 1grams (PDF)

S1 Table. S1 Table (a) and (b) show the number of 1-grams of each type gathered from the *All English* and *British-English* Googled Ngram Databases for this study, as well as the possible associated *Lexical Landscapes* that could be constructed from each set. S1 Table (c) and (d) likewise show the number of available words and the total possible landscapes that could be constructed from the reduced data sets used within this investigation. The reduced data sets were constructed using cutoffs in order to eliminate words with excessively low usage counts and avoid the presence of acronyms within *Lexical Landscapes*. (PDF)

S1 Fig. Dis-aggregated epistasis across Lexical Landscapes. In the main text Fig 6 demonstrated the aggregated effects, the coefficients corresponding to 0th—4th order effects. S2(a) and S2(b) Fig show how individual epistatic effects can vary over time, or more generally, across context. Each line represents a particular epistatic effect, and the lines are grouped by order and demarcated by color theme and marker. For instance, the collection of light red to dark red lines shows the *first* order epistatic effects: 0001, 0010, 0100, 1000.

(TIF)

S2 Fig. Aggregated British-English higher-order epistatic effects. Fitness graphs for the POEM → LAST landscape. As described in this manuscript, the aggregated epistatic effects combine the averages of individual interactions and organize them by their order. S3 Fig shows these effects for the British-English CARS → SOME and POEM → LAST landscapes.

(TIF)

S3 Fig. Dis-aggregated British-English higher-order epistatic effects. As discussed in several places throughout this manuscript, dis-aggregated graphs represent how individual epistatic terms interact across contexts. These graphs represent those effects for the CARS → SOME and POEM → LAST landscapes in British-English.

(TIF)

S4 Fig. British-English within-path competition. Here we present within-path competition in the British English Lexical Landscapes subset.

(TIF)

S5 Fig. Fitness graphs for 3 letter (HAY → FIT) and 5 letter (ADDED → VIRUS) Lexical Landscapes. Fitness graphs for 3 letter (HAY → FIT) and 5 letter (ADDED → VIRUS) Lexical Landscapes. S6(a) and S6(b) Fig are visualizations of the fitness landscapes for three and five letters. The color indicates the fitness: darker the blue, higher the fitness. Edges in the graph are weighted by the difference (in absolute value) between the fitness values of the two adjacent nodes and are emboldened in a proportional way to show the weight.

(TIF)

S6 Fig. Within-path competition (C_w) for 3 and 5 world landscapes. S7(a) and S7(b) Fig demonstrate that the (C_w) can vary across environment (as it did for the 4 letter 1-grams discussed in the main texts). Red dots indicate the years where the specified path—shown to the right of each figure—represented an *uphill* path in the fitness landscape, that is, a path where each successive word in the path has a higher fitness than the previous one.

(TIF)

S7 Fig. Epistasis across environment. S7 Fig (a) and (b) show how epistatic effects for example three and five letter word landscapes can vary across environment, or in this case across time. We have chosen a three letter landscape (BAD to DRY) with mild fluctuations over time to contrast it with the relatively large fluctuations in the five-letter landscape (SHARP to ATONE).

(TIF)

Acknowledgments

The authors would like to thank D. Weinreich, A. Alexander and S.J. Gates for helpful discussion and input on the project.

Author Contributions

Conceptualization: C. Brandon Ogbunugafor.

Data curation: Victor A. Meszaros, Miles D. Miller-Dickson.

Formal analysis: Victor A. Meszaros, Miles D. Miller-Dickson, C. Brandon Ogbunugafor.

Funding acquisition: Miles D. Miller-Dickson, C. Brandon Ogbunugafor.

Investigation: Victor A. Meszaros, Miles D. Miller-Dickson, C. Brandon Ogbunugafor.

Methodology: Victor A. Meszaros, Miles D. Miller-Dickson.

Project administration: C. Brandon Ogbunugafor.

Resources: Victor A. Meszaros, Miles D. Miller-Dickson.

Software: Victor A. Meszaros, Miles D. Miller-Dickson.

Supervision: C. Brandon Ogbunugafor.

Validation: Victor A. Meszaros.

Visualization: Victor A. Meszaros, Miles D. Miller-Dickson.

Writing – original draft: Victor A. Meszaros, Miles D. Miller-Dickson, C. Brandon Ogbunugafor.

Writing – review & editing: Victor A. Meszaros, Miles D. Miller-Dickson, C. Brandon Ogbunugafor.

References

1. Provine WB. The origins of theoretical population genetics: With a new afterward. University of Chicago Press. 2001.
2. Nowak MA. Evolutionary dynamics. Harvard University Press. 2006.
3. Wright S. The roles of mutation, inbreeding, crossbreeding, and selection in evolution. Proceedings of the Sixth International Congress of Genetics. 1932; 1:356–366.
4. Gavrilets S. Fitness landscapes and the origin of species (MPB-41). Princeton University Press. 2004; 41.
5. De Visser JAG, and Kurg J. Empirical fitness landscapes and the predictability of evolution. Nature Review Genetics. 2014; 15:480 <https://doi.org/10.1038/nrg3744>
6. Weinreich DM, Delaney NF, DePristo MA, and Ha DL. Darwinian evolution can follow only very few mutational paths to fitter proteins. Science. 2006; 312:111–114 <https://doi.org/10.1126/science.1123539> PMID: 16601193
7. Weinreich DM, Lan Y, Wylie SC, and Heckendorn RB. Should evolutionary geneticists worry about higher-order epistasis?. Current opinion in genetics & development, Elsevier. 2013; 23(6):700–707. <https://doi.org/10.1016/j.gde.2013.10.007>
8. Ogbunugafor CB and Eppstein MJ. Competition along trajectories governs adaptation rates towards antimicrobial resistance. Nature Ecology and Evolution. 2017; 1(7). <https://doi.org/10.1038/s41559-016-0064> PMID: 28812621
9. Palmer AC and Kishon R. Understanding, predicting and manipulating the genotypic evolution of antibiotic resistance. Nature Review Genetics 2013; 14(243).
10. Nichol D, Jeavons P, Fletcher AG, Bonomo RA, Maini PK, Paul JL, et al. Steering evolution with sequential therapy to prevent the emergence of bacterial antibiotic resistance. PLoS Computational Biology. 2015; 11(e10044935). <https://doi.org/10.1371/journal.pcbi.1004493>
11. Palmer AC, Toprak E, Baym M, Kim S, Veres A, Bershtein S, et al. Delayed commitment to evolutionary fate in antibiotic resistance fitness landscapes?. Nature Communication. 2015; 6 (7385). <https://doi.org/10.1038/ncomms8385>

12. Aguilar-Rodríguez J, Payne JL, and Wagner A A thousand empirical adaptive landscapes and their navigability. *Nature Ecology and Evolution*. 2017; 1(45). <https://doi.org/10.1038/s41559-016-0045-28812623> PMID: 28812623
13. Weinreich DM, Lan Y, Jaffe J, and Heckendo RB. The influence of higher-order epistasis on biological fitness landscape topography. *Journal of Statistical Physics*, Elsevier. 2018; 172:208–225. <https://doi.org/10.1007/s10955-018-1975-3>
14. Baeza-Centurion P, Miñana B, Schmiedel JM, Valcárcel J, and Lehner B. Combinatorial genetics reveals a scaling law for the effects of mutations on splicing. *Cell*. 2019. <https://doi.org/10.1016/j.cell.2018.12.010> PMID: 30661752
15. Wilke CO, Wang JL, Ofria C, Lenski RE, and Adami C. Evolution of digital organisms at high mutation rates leads to survival of the flattest. *Nature* 2001; 412:331. <https://doi.org/10.1038/35085569> PMID: 11460163
16. Brouillet S, Annoni H, Ferretti L, and Achaz G. Magellan: a tool to explore small fitness landscapes. *BioRxiv*. 2015;031583.
17. Smith JM. Natural selection and the concept of a protein space. *Nature*. 1970; 225:563. <https://doi.org/10.1038/225563a0> PMID: 5411867
18. Ogbunugafor CB and Hartl DL. A new take on john maynard smith's concept of protein space for understanding molecular evolution. *PLoS Computational Biology*. 2016; 12(e1005046).
19. Michel J-B, Shen YK, Aiden AP, Veres A, Gray MK, Pickett JP, et al. Quantitative analysis of culture using millions of digitized books. *Science*. 2011; 331:176–182. <https://doi.org/10.1126/science.1199644> PMID: 21163965
20. Lin Y, Michel JB, Aiden EL, Orwant J, Brockman W, and Petrov S. Syntactic annotations for the google books ngram corpus. *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*. 2012:169–174.
21. Eppstein MJ. DARPS. <http://www.cs.uvm.edu/meppstei/DARPS>. 2016.
22. Poelwijk FJ, Krishna V, and Ranganathan R. The Context-Dependence of Mutations: A Linkage of Formalisms. *PLoS Computational Biology*. 2016; 12(e1004771) <https://doi.org/10.1371/journal.pcbi.1004771> PMID: 27337695
23. Crawford L, Zeng P, Mukherjee S, and Zhou X. Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. *PLoS Genetics*. 2017; 13(e1006869). <https://doi.org/10.1371/journal.pgen.1006869> PMID: 28746338
24. Otwinowski J, McCandlish DM, and Plotkin JB. Inferring the shape of global epistasis. *Proceedings National Academy of Science*. 2018; 115:E7550–E7558. <https://doi.org/10.1073/pnas.1804015115>
25. Lozovsky ER, Chookajorn T, Brown KM, Imwong M, Shaw PJ, Kamchonwongpaisan S, et al. Stepwise acquisition of pyrimethamine resistance in the malaria parasite. *Proceedings National Academy of Science*. 2009; 106:12025–12030. <https://doi.org/10.1073/pnas.0905922106>
26. Ogbunugafor CB, Wylie CS, Diakite I, Weinreich DM, and Hartl DL. Adaptive landscape by environment interactions dictate evolutionary dynamics in models of drug resistance. *PLoS Computational Biology*. 2016; 12(e1004710). <https://doi.org/10.1371/journal.pcbi.1004710> PMID: 26808374
27. Iwasa Y, Michor F, and Nowak MA. Stochastic tunnels in evolutionary dynamics. *Genetics*. 2004; 166:1571–1579. <https://doi.org/10.1534/genetics.166.3.1571> PMID: 15082570
28. Weinreich DM, Lan Y, Wylie SC, and Heckendorn RB. Should evolutionary geneticists worry about higher-order epistasis?. *Current opinion in genetics & development*, Elsevier. 2013; 23(6):700–707. <https://doi.org/10.1016/j.gde.2013.10.007>
29. Flynn KM, Cooper TF, Moore FB, and Cooper VS. The environment affects epistatic interactions to alter topology of an empirical fitness landscape. *PLoS Genetics*. 2013; 9(e1003426). <https://doi.org/10.1371/journal.pgen.1003426> PMID: 23593024
30. Sackton TB and Hartl DL. Genotypic context and epistasis in individuals and populations. *Cell*. 2016; 166:279–287. <https://doi.org/10.1016/j.cell.2016.06.047> PMID: 27419868
31. Guerrero RF, Scarpino SV, Rodrigues JV, Hartl DL, and Ogbunugafor CB. Proteostasis environment shapes higher-order epistasis operating on antibiotic resistance. *Genetics*. 2019;302138. <https://doi.org/10.1534/genetics.119.302138> PMID: 31015194