# Population Diversity Leads to Short Running Times of Lexicase Selection

Thomas Helmuth[1], Johannes Lengler[2], and William La Cava[3*]

[1] Hamilton College, Clinton, New York, USA
thelmuth@hamilton.edu
[2] ETH Zürich, Zürich, Switzerland
[3] Boston Children's Hospital, Harvard Medical School, Boston, MA, USA
[*]corresponding author: william.lacava@childrens.harvard.edu

**Abstract.** In this paper we investigate why the running time of lexicase parent selection is empirically much lower than its worst-case bound of $O(N \cdot C)$. We define a measure of population diversity and prove that high diversity leads to low running times $O(N + C)$ of lexicase selection. We then show empirically that genetic programming populations evolved under lexicase selection are diverse for several program synthesis problems, and explore the resulting differences in running time bounds.

**Keywords:** lexicase selection · population diversity · running time analysis.

## 1 Introduction

Semantic selection methods have been of increased interest as of late in the evolutionary computation community [16,22] due to the observed improvements over more traditional selection methods (such as tournament selection) that only consider individual behavior in aggregate. One such method is lexicase selection [19,10], a parent selection method originally proposed for genetic programming. Since then, the original algorithm and its variants have found success in different domains, including program synthesis [8], symbolic regression [18,13], evolutionary robotics [17], and learning classifier systems [1].

Although an active research community has illuminated many aspects of lexicase selection's behavior via experimental analyses [5,4,17,7], theoretical analyses of lexicase selection have been slower to develop. Previous theoretical work has looked at the probability of selection under lexicase, and also made connections between lexicase selection and Pareto optimization [12]. A study focusing on ecological theory provided insights into the efficacy of lexicase selection [3]. Additionally, the running time of a simple hill climbing algorithm utilizing lexicase selection has been analyzed for the bi-objective leading ones trailing zeroes benchmark problem [11]. However, the recursive nature of lexicase selection, and its step-wise dependence on the behavior of subsets of the population, make it difficult to analyze.

We focus this paper on a particular gap in the theory of lexicase selection, which is an understanding of its running time. Although the worst-case complexity is known to be $O(N \cdot C)$, where $N$ is the population size and $C$ is the set of training cases, empirical data suggest the worst-case condition is extremely rare in practice [12]. Our goal is to explain this discrepancy through a combination of theory and experiment.

### 1.1   Our contributions

We find that the observed running time of lexicase selection can be explained with *population diversity*, by which we mean the phenotypic/behavioral diversity of individuals in a population. Our contributions are threefold:

1. We introduce a new way of measuring population diversity, Definitions 2 and 3, which we call $\varepsilon$-*Cluster Similarity*, or $\varepsilon$-*Similarity* for short. Here, for different values of the parameter $\varepsilon$, we obtain a measure of how similar the population is, where small $\varepsilon$-Cluster Similarity corresponds to high diversity. As we show, this measure is not directly tied to other measures of diversity like the average phenotypical distance (Section 2.2) or the mean of the behavioral covariance matrix (Fig. 3).
2. We prove mathematically that lexicase selection is fast when applied to populations which are diverse. More precisely, we show that with low $\varepsilon$-Cluster Similarity the expected running time of lexicase selection drops from $O(N \cdot C)$ to $O(N + C)$, where the hidden constants depend on the parameter $\varepsilon$ and on the quantity $k$ that measures $\varepsilon$-Cluster Similarity.
3. Finally, we show empirically for several program synthesis problems [8] that genetic programming populations are indeed diverse in our sense (have low $\varepsilon$-Similarity). We investigate which parameter $\varepsilon$ gives the best running time guarantees for lexicase selection, and we find that the running time guarantees are substantially better than the trivial running time bound of $N \cdot C$.

Our findings apply to both discrete and continuous problems and population behaviors. Although we restrict our analysis to vanilla lexicase selection, we note the results generalize to other variants, including $\epsilon$-lexicase selection [14][4] and down-sampled lexicase selection [9].

## 2   Preliminaries

### 2.1   Lexicase Selection

Lexicase selection is used to select a parent for reproduction in a given population. Unlike many common parent selection methods, lexicase selection does not aggregate an individual's performance into a single fitness value. Instead, it

---

[4] Our results hold for the original variant, later dubbed "static" $\varepsilon$-lexicase selection [12].

---

**Algorithm 1 Lexicase Selection** applied to a population $\mathcal{N}$ without dupli-cates, with discrete loss/error $L(n, c)$ on training cases $c \in \mathcal{C}$ and individual $n \in \mathcal{N}$. Returns an individual selected to be a parent. $\mathcal{N}_t$ is the remaining candidate pool at step $t$, $\mathcal{C}'$ is the set of remaining training cases.

---

**LEX**($\mathcal{N}$, $\mathcal{C}$, $L$):
    $\mathcal{C}' \leftarrow \mathcal{C}$; $t \leftarrow 0$; $\mathcal{N}_0 \leftarrow \mathcal{N}$;
    **while** $|\mathcal{N}_t| > 1$:
        $c \leftarrow$ random choice from $\mathcal{C}'$
        $\ell^* \leftarrow \min\{L(n, c) \mid n \in \mathcal{N}_t\}$
        $\mathcal{N}_{t+1} \leftarrow \{n \in \mathcal{N}_t \mid L(n, c) = \ell^*\}$
        $\mathcal{C}' \leftarrow \mathcal{C}' \setminus \{c\}$
        $t \leftarrow t + 1$
    **return** unique element from $\mathcal{N}_t$

---

considers the loss (errors) on different training *cases* (a.k.a. samples/examples) independently, never comparing (even indirectly) the results on one training case with those on another.

Lexicase selection begins by putting the entire population into a candidate pool. As a preprocessing step, all phenotypical *duplicates* are removed from the pool, i.e., if several individuals give the same loss on all training cases, all but one are removed in the following filtering steps. Then lexicase selection repeatedly selects a new training case $t$ at random, and removes all individuals from the current candidate pool that do not achieve the best loss on case $t$ within the current pool. This process is repeated until the candidate pool contains only a single individual. If the remaining individual has phenotypical duplicates, the selected parent is taken at random from among these behavioral clones. We formalize the algorithm in Alg. 1.

We remark that the process is guaranteed to end up with a candidate pool of size one: whenever the candidate pool contains at least two individuals, they are not duplicates due to preprocessing. Hence, they differ on at least one training case $c$, and one of them is filtered out when $c$ is considered. So after all training cases have been processed, it is not possible that the candidate pool contains more than one individual.

Note that the described procedure selects a single individual from the popu-lation. In order to gather enough parents for the next generation, it is typically performed $O(N)$ times, where $N$ is the population size. An exception is the pre-processing step that only needs to be performed once each generation. Moreover, finding duplicates can be efficiently implemented via a hash map. Thus prepro-cessing is usually not the bottleneck of the procedure, and we will focus in this paper on the remaining part: the repeated reduction of the selection pool via random training cases. To exclude the effect of preprocessing, we will assume that the initial population is already free of duplicates.

In case of real-valued (non-discrete) losses, one typically uses a variant known as $\epsilon$-lexicase selection [14]. (Note this use of $\epsilon$ is distinct from that used in Sec-tion 2.2). In the original algorithm, later dubbed "static" $\epsilon$-lexicase selection [12],

phenotypic behaviors are binarized prior to lexicase selection, such that individuals within $\epsilon$ of the population-wide minimal loss on $c$ have an error of 0, and otherwise an error of 1. Our results extend naturally to this version of $\epsilon$-lexicase selection.

In contrast, in the "dynamic" and "semi-dynamic" variants of $\epsilon$-lexicase selection, lexicase selection removes those individuals whose loss is larger than $\ell^* + \epsilon$, where $\ell^*$ is the minimal loss in the *current candidate pool* [12]. Our results may extend to these scenarios, but the framework becomes more complicated. Here, it is no longer possible to separate the preprocessing step (i.e., de-duplication) from the actual selection mechanism. Of course, it is still possible to define two individuals $n_1, n_2 \in \mathcal{N}$ as *duplicates* if they differ by at most $\delta$ on all training cases. But this is no longer a transitive relation, i.e., it may happen that $n_1$ and $n_2$ are duplicates, $n_2$ and $n_3$ are duplicates, but $n_1$ and $n_3$ are not duplicates. For these reasons, it is necessary to handle duplicates indirectly during the execution of the algorithm. To avoid these complications, we only present Alg. 1 in the case of discrete losses and without duplicates, but we do include the case of real-valued losses in our analysis.

The worst-case running time of lexicase selection is $O(N \cdot C)$, where $N := |\mathcal{N}|$ is the population size and $C := |\mathcal{C}|$ is the number of training cases. The problem is that in an iteration of the while-loop, it may happen that $\mathcal{N}_{t+1} = \mathcal{N}_t$, i.e., that the candidate pool does not shrink. This is more likely for binary losses. Then, it may happen that $\mathcal{N}_t = \mathcal{N}$ for many iterations of the while-loop, and then computing $\ell^*$ and $\mathcal{N}_{t+1}$ needs time $O(N)$. Since the while-loop is executed up to $C$ times, this leads to the worst-case runtime $O(N \cdot C)$. Recall that we usually want to run the procedure $O(N)$ times to select all parents for the next generation, which then takes time $O(N^2 \cdot C)$.

One might hope that the expected runtime is much better than the worst-case runtime. This is the case for many classical algorithms like quicksort, but for populations with an unfavorable loss profile, it is not the case here. Consider a population of individuals which have the same losses on all training cases except a single case $c$, and on case $c$ they all have different losses. Then the candidate pool does not shrink before this case $c$ is found, and finding this case needs $C/2$ iterations in expectation. Thus the expected runtime is still of order $O(N \cdot C)$.[5]

So in order to give better bounds on the running time of lexicase selection, it is required to have some understanding of the involved populations. This is precisely the contribution of this paper: we define a notion of diversity that provably leads to a small running time of lexicase selection, and we empirically

---

[5] This worst-case example does not hold if the losses are binary, but even that does not help much. It is possible to construct a population of $N$ individuals without duplicates that differ only on $\log_2 N$ binary training cases, and are identical on all other training cases. In this situation, the candidate pool does not shrink before at least one of those training cases is found, and in expectation this takes $C/\log_2 N$ iterations. Thus the expected runtime in this situation is at least $O(N \cdot C/\log N)$, which is not much better than $O(N \cdot C)$.

show that populations in genetic programming are diverse with respect to this measure.

## 2.2  $\varepsilon$-Cluster Similarity

We now come to our first main contribution, a new way of measuring diversity. The measure is in phenotype space, so it measures for a training set $\mathcal{C}$ how similar the individuals perform on this training set. We first introduce a useful notion, which is the *phenotypical distance* of two individuals.

**Definition 1 (Phenotypical Distance).** *Consider two individuals $m, n$ that are compared on a set $\mathcal{C}$ of training cases. The* phenotypical distance *between $m$ and $n$ is the number of training cases in which $m$ and $n$ have different losses.*

*If the losses are real-valued, then for $\delta > 0$, the* phenotypical $\delta$-distance *between $m$ and $n$ is the number of training cases in which the losses of $m$ and $n$ differ by at least $\delta$.*

**Definition 2 ($\varepsilon$-Cluster Similarity).**  *Let $\mathcal{N}$ be a population of individuals, and $\mathcal{C}$ be a set of training cases with discrete losses, for example binary losses. Let $\varepsilon \in [0, 1]$. Then the $\varepsilon$-Cluster Similarity is defined to be the minimal $k \geq 2$ such that among every set of $k$ different individuals in $\mathcal{N}$, there are at least two individuals $m, n \in \mathcal{N}$ with phenotypical distance at least $\varepsilon|\mathcal{C}|$.*

*If instead $\mathcal{C}$ is a training set with real-valued losses, then let $\varepsilon \in [0, 1]$ and $\delta > 0$. Then the $\varepsilon$-Cluster Similarity for $\delta$-distance is defined as the minimal $k \geq 2$ such that among every set of $k$ different individuals in $\mathcal{N}$, there are at least two individuals $m, n \in \mathcal{N}$ with phenotypical $\delta$-distance at least $\varepsilon|\mathcal{C}|$.*

A few remarks are in order to understand the definition better. Firstly, the $\varepsilon$-Cluster Similarity $k$ is a *decreasing* measure of diversity, i.e., less similarity means more diversity and vice versa. Moreover, the value $k$ is increasing in $\varepsilon$: we are only satisfied with two individuals of distance at least $\varepsilon|\mathcal{C}|$, which is harder to achieve for larger values of $\varepsilon$. Therefore, we may need a larger set to ensure that it contains a pair of individuals with such a large distance. In other words, a larger value of $\varepsilon$ means that we are more restrictive in counting individuals as different, which yields larger value of $k$: the population is more similar with respect to a more restrictive measure of difference. There is an important tradeoff between $\varepsilon$ and $k$: larger values of $\varepsilon$ (which are desirable in terms of diversity, since we search for individuals with larger distances) lead to larger values of $k$ (which is undesirable since we only find such individuals in larger sets).

Second, having small $\varepsilon$-Cluster Similarity is a rather weak notion of diversity: it does not require that *all* pairs of individuals are different from each other. For example, if the population consists of clusters of $k - 1$ individuals which are pairwise very similar, then the $\varepsilon$-Cluster Similarity is $k$ as long as the clusters have distances at least $\varepsilon|\mathcal{C}|$ from each other. We just forbid that there is a cluster of size $k$ such that *every* pair of individuals in the cluster has small distance.

On the other hand, the $\varepsilon$-Cluster Similarity may be a finer measure than, say, the average phenotypical distance in the population. For example, consider

a population that consists only of two clusters of almost identical individuals, but the clusters are in opposite corners of the phenotype space, i.e., they differ on almost all training cases. Then the average phenotypical distance is extremely large, $\approx |\mathcal{C}|/2$, which would suggest high diversity. But even for absurdly high $\varepsilon = 0.9$, we would find $k = |\mathcal{N}|/2 + 1$, i.e., a very low diversity according to our definition. It is not hard to see that in this example the expected running time of lexicase selection is $\Omega(|\mathcal{N}| \cdot |\mathcal{C}|)$: in the first step one of the clusters will be removed completely, but afterwards it is very hard to make any further progress. Hence, this example shows that *average phenotypical distance does not predict the running time of lexicase selection well*: even though the example has large average phenotypical distance ("large diversity" in that sense), the running time of lexicase selection is very high. The main theoretical insight of this paper is that this discrepancy can never happen with $\varepsilon$-Cluster Similarity. Whenever $\varepsilon$-Cluster Similarity is low (large diversity), then the expected running time of lexicase selection is small.

To give the reader another angle to grasp the definition of $\varepsilon$-Cluster Similarity, we give a second, equivalent definition in terms of graph theory.

**Definition 3 ($\varepsilon$-Cluster Similarity, Equivalent Definition).** *Let $\mathcal{N}$ be a population of search points, and $\mathcal{C}$ be a set of training cases with discrete losses. Let $\varepsilon \in [0,1]$. We define a graph $G = (V, E)$ as follows. The vertex set $V := \mathcal{N}$ is identical with the population. Between any two vertices $m, n \in \mathcal{N}$, we draw an edge if and only if the individuals $m$ and $n$ have the same loss in more than $(1 - \varepsilon)|\mathcal{C}|$ training cases. Then the $\varepsilon$-Cluster Similarity is $k := \alpha + 1$, where $\alpha$ is the clique number of $G$, i.e, $\alpha$ is the size of the largest clique of the graph $G$.*

*If $\mathcal{C}$ is a training set with real-valued losses and $\delta > 0$ a parameter, then we use the same vertex set for $G = G(\delta)$, but we draw an edge between $m$ and $n$ if and only if the losses of $m$ and $n$ differ by at most $\delta$ for more than $(1 - \varepsilon)|\mathcal{C}|$ training cases. Then the $\varepsilon$-Cluster Similarity for $\delta$-distance is again $k := \alpha + 1$, where $\alpha$ is the clique number of $G(\delta)$.*

## 3    Theoretical Result: Low $\varepsilon$-Cluster Similarity leads to small running times

In this section, we prove mathematically, that a high diversity (i.e., a small $\varepsilon$-Similarity) leads to a small expected running time for lexicase selection.

### 3.1    Preliminaries

To proof our main theoretical result, we will use the following theorem, known as Multiplicative Drift Theorem [2,15], which is a standard tool in the theory of evolutionary computation.

**Theorem 1 (Multiplicative Drift).** *Let $(X_t)_{t \geq 0}$ be a sequence of non-negative random variables with a finite state space $\mathcal{S} \subseteq \mathbb{R}_0^+$ such that $0 \in \mathcal{S}$. Let $s_{\min} :=$*

$\min(\mathcal{S} \setminus \{0\})$, *let* $T := \inf\{t \geq 0 \mid X_t = 0\}$, *and for* $t \geq 0$ *and* $s \in \mathcal{S}$ *let* $\Delta_t(s) := E[X_t - X_{t+1} \mid X_t = s]$. *Suppose there exists* $\delta > 0$ *such that for all* $s \in \mathcal{S} \setminus \{0\}$ *and all* $t \geq 0$ *the drift is*

$$\Delta_t(s) \geq \delta s. \tag{1}$$

*Then*

$$E[T] \leq \frac{1 + E[\ln(X_0/s_{\min})]}{\delta}. \tag{2}$$

Now we can give our theoretical results. Note that the following theorems refer to a single execution of lexicase selection, i.e., $M$ refers to the complexity of finding a single parent via lexicase selection. The following theorem says that the running time is low, $O(|\mathcal{N}| + |\mathcal{C}|)$ if the population has large $\varepsilon$-Cluster Similarity. As common in theoretical running time analysis, we give the running time in terms of *evaluations*, where an evaluation is an execution (or lookup) of $L(n, c)$ for an individual $n$ and a training case $c$. The running time is proportional to the number of evaluations.

**Theorem 2.** *Let* $0 < \varepsilon < 1$. *Consider lexicase selection on a population* $\mathcal{N}$ *without duplicates and with* $\varepsilon$-*Cluster Similarity of* $k \in \mathbb{N}$. *Let* $M$ *be the number of evaluations until the population pool is reduced to size* $1$. *Then*

$$E[M] \leq \frac{4|\mathcal{N}|}{\varepsilon} + 2k|\mathcal{C}|.$$

*Proof.* Consider any two individuals $m, n \in \mathcal{N}$. Assume that both $m, n$ are still in the candidate pool after some selection steps (i.e., after some iterations of the while-loop have been processed). Then $m$ and $n$ can not differ in any of the processed cases $\mathcal{C} \setminus \mathcal{C}'$, because otherwise one of them would have been removed from the population. Therefore, the $\varepsilon|\mathcal{C}|$ cases in which $m$ and $n$ differ are all still contained in $\mathcal{C}'$. In particular, if we choose a new case from $\mathcal{C}'$ at random, then the probability that $m$ and $n$ differ in this case is at least $\varepsilon|\mathcal{C}|/|\mathcal{C}'| \geq \varepsilon$. Note that this holds throughout the algorithm and for any two individuals $m, n$ that are still candidates.

Now we turn to the computation. Let $X_t$ be the number of remaining individuals after $t$ executions of the while-loop. We define $Y_t := X_t$ if $X_t \geq 2k$ and $Y_t := 0$ if $X_t < 2k$. Let $T'$ be the first point in time when $Y_{T'} = 0$ (and thus, $X_{T'} < 2k$).

If $X_t \geq 2k$, then we split the population before the $t + 1$-st step into pairs as follows. Since $X_t \geq k$, there are at least two individuals which differ in at least $\varepsilon|\mathcal{C}|$ cases, so we pick two such individuals and pair them up. We can iterate this until there are less than $k$ unpaired individuals left. Therefore, we are able to pair up at least $X_t - (k - 1) > X_t - X_t/2 = X_t/2$ individuals, forming at least $X_t/4$ pairs. For each pair, there is a chance of $\varepsilon$ that the two differ in the case of the $t + 1$-st step, in which case at least one of them is eliminated. Hence, for every $x \geq 2k$,

$$E[X_{t+1} \mid X_t = x] \leq x - \tfrac{\varepsilon x}{4} = x(1 - \tfrac{\varepsilon}{4}).$$

Now let us assume that $Y_t = y > 0$ (and thus $y \geq 2k$). Since $Y_{t+1} \leq X_{t+1}$ by definition, we obtain

$$E[Y_{t+1} \mid Y_t = y] \leq E[X_{t+1} \mid X_t = y] \leq y(1 - \tfrac{\varepsilon}{4}). \qquad (3)$$

The advantage of $Y_t$ is that the above bound holds for all $y \geq 0$ (it is trivial for $y = 0$), whereas the corresponding bound for $X_t$ may not hold for $0 < x < 2k$.

Now we bound $M$ by splitting it into the running time $M_1$ before step $T'$, and the running time $M_2$ after and including step $T'$. For $M_1$, we proceed as follows.

$$M_1 = \sum_{t=0}^{T'-1} X_t = \sum_{t=0}^{T'-1} Y_t = \sum_{t=0}^{\infty} Y_t,$$

because $Y_t = 0$ for $t \geq T'$. Applying (3) iteratively to $Y_t$, we obtain

$$E[Y_t] \leq (1 - \tfrac{\varepsilon}{4})^t Y_0,$$

where $Y_0 = |\mathcal{N}|$. Plugging this in, we get

$$E[M_1] = \sum_{t=0}^{\infty} E[Y_t] \leq \sum_{t=0}^{\infty} (1 - \tfrac{\varepsilon}{4})^t (|\mathcal{N}| - 1) \leq \tfrac{4|\mathcal{N}|}{\varepsilon},$$

where in the last step we have used the formula $\sum_{t=0}^{\infty} q^t = 1/(1-q)$ for geometric series with $q = 1 - \varepsilon/4$. It remains to bound $M_2$, and we use a simple bound. Since every case occurs at most once, and since the population size is at most $2k$, we have deterministically $M_2 \leq 2k \cdot |\mathcal{C}|$. $\qquad \square$

## 4   Empirical Evaluation in Program Synthesis

We evaluated the theoretical bounds given by Theorem 2 on examples using genetic programming to solve program synthesis benchmark problems. The purpose of this evaluation is to 1) find out how diverse the populations are according to $\varepsilon$-Cluster Similarity; 2) measure the extent to which the new bounds shrink our estimates of the running time of lexicase selection, relative to the known worst-case bounds; 3) evaluate the sensitivity of $\varepsilon$-Cluster Similarity to the parameter, $\varepsilon$, across several problems; and 4) determine how $\varepsilon$-Cluster Similarity compares to a more standard diversity metric in real data.

### 4.1   Experimental Setup

We investigate these aims using 8 program synthesis problems taken from the General Program Synthesis Benchmark Suite [8]. These problems require solution programs to manipulate multiple data types and exhibit different control structures, similar to the types of programs we expect humans to write. Among the 8 problems there are 5 different expected output types (Boolean, integer, float, vector of integers, and string), allowing us to test against multiple data types. In particular, we note that two problems (compare-string-lengths and

mirror-image) have Boolean outputs, which we expect to have higher $\varepsilon$-Cluster Similarity values due to having fewer possible output values.

Our experiments were conducted using PushGP, which evolves programs in the Push programming language [21,20]. PushGP is expressive in the data types and control structures it allows, and has been used previously with these problems [8]. We use the Clojush, the Clojure implementation of Push, in our experiments.[6] Each run evolves a population of 1000 individuals for a maximum of 300 generations using lexicase selection and UMAD mutation (without crossover) [6]. We conduct 100 runs of genetic programming per problem.

For each problem, trial, generation, and selection event, we calculated the $\varepsilon$-Cluster Similarity for $\varepsilon \in [0.05, 0.6]$ in increments of 0.05. Using these values, we calculated the bound on the expected running time of lexicase selection according to Theorem 2. For comparison, we also calculated 1) the worst-case complexity of lexicase selection at those operating points and 2) the average pair-wise covariance of the population error vectors. We calculated worst-case running time as $N \cdot C$, neglecting constants, in order to make our comparison to the new running time calculation conservative.

### 4.2   Results

Fig. 1 visualizes the new running time bound as a fraction of the bound given by the worst-case complexity, $N \cdot C$. Across problems, the running time bound given by Theorem 2 ranges from approximately 10-70% of the worst-case complexity bound, indicating much lower expected running times. On average over all problems, the bound given by Theorem 2 is 24.7% of the worst-case bound on running time.

Fig. 2 shows the components of Theorem 2 as a function of $\varepsilon$, as well as the total expected running time bound. For small values of $\varepsilon$, the $4|\mathcal{N}|/\varepsilon$ term dominates, whereas for larger values of $\varepsilon$, the $2k|\mathcal{C}|$ term dominates. The observed behavior agrees with our intuition, since larger values of $\varepsilon$ lead to larger values of $k$. The value of $\varepsilon$ corresponding to the lowest bound on running time varies by problem, with an average value of 0.29.

Fig. 3 compares the new diversity metric (Definition 3) to a more typical definition of behavioral diversity: the mean of the covariance matrix given by population errors. In general, we observe that $\varepsilon$-Cluster Similarity does not correlate strongly with mean covariance, suggesting that it does indeed measure a different aspect of phenotypic diversity as suggested in Section 2.2.

## 5   Discussion

We see in Fig. 1 that the new running time bound is below the old bound, and sometimes substantially lower by a factor of $5 - 10$. Thus a substantial part of the discrepancy between the (old) worst-case running time bound and
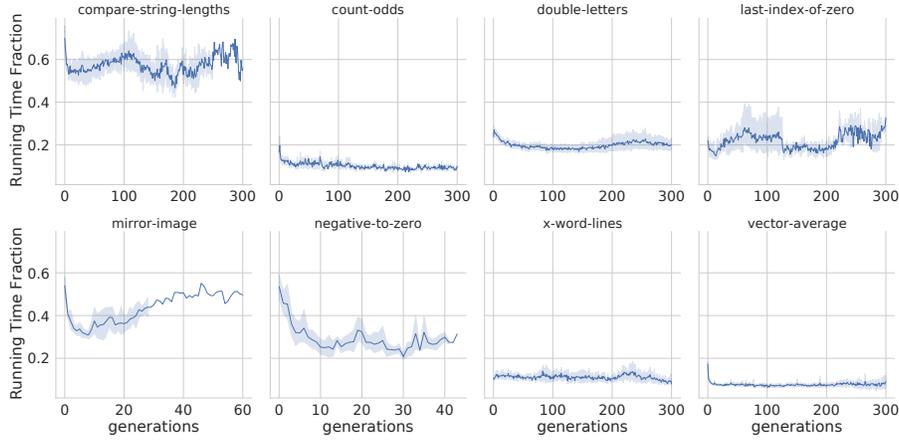
---

[6] https://github.com/lspector/Clojush

**Fig. 1.** New running time bound divided by the previously known worst-case bound, as a function of evolutionary generations. The $y$-axis shows the ratio of both bounds using measurements of relevant parameters. The filled region represents confidence interval of the estimates over all trials.
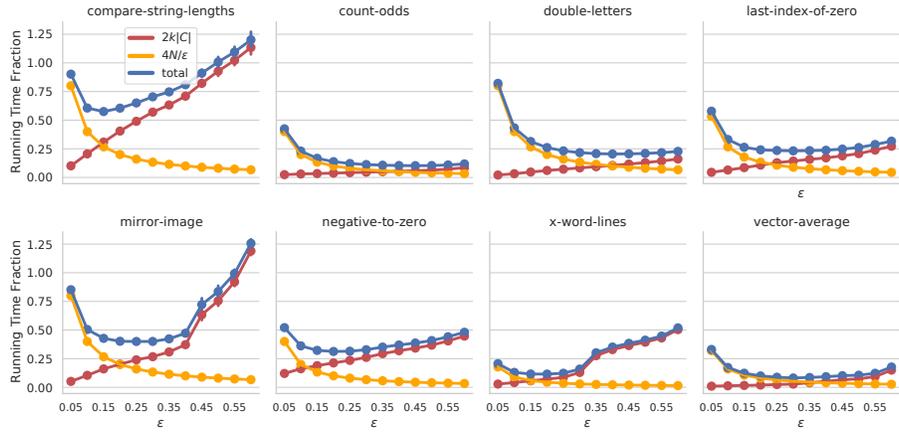


**Fig. 2.** The ratio of new running time bound and previous known worst-case bound as a function of $\varepsilon$. Optimal values vary by problem but we note flat regions for many problems suggesting a broad range of possible $\varepsilon$ values that give similar running time bounds.

the empirically observed fast running times can be explained by the fact that populations in real-world data are diverse according to our new measure. Since the theoretical analysis is still a worst-case analysis (over all populations), we
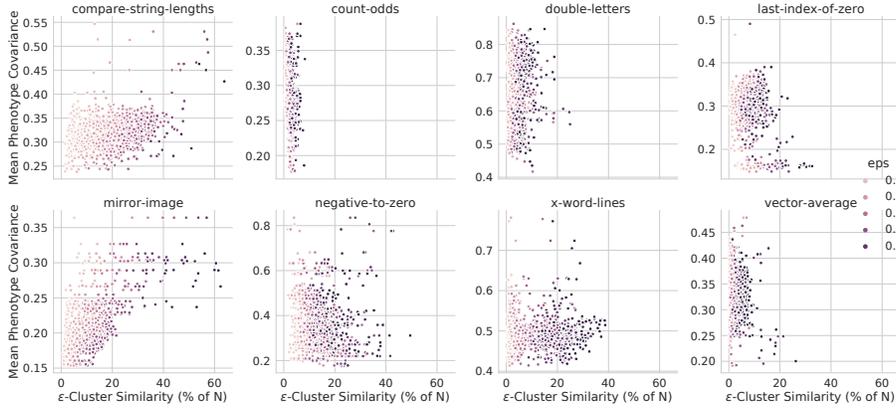
**Fig. 3.** A comparison of $\varepsilon$-Cluster Similarity (x-axis) and the mean of the covariance matrix of population error (y-axis), colored by $\varepsilon$. $\varepsilon$-Cluster Similarity ($k$) is plotted as a percent of the population size, $N$. In most cases (6/8), we observe little relation between the two measures, suggesting $\varepsilon$-Cluster Similarity is indeed measuring a distinct aspect of population diversity, one that is particularly relevant to the running time of lexicase selection.

do not expect that the new bound can explain the whole gap in all situations, but it does explain a substantial factor.

In Fig. 2, we investigate which choice of $\varepsilon$ gives the best running time bound. Note that $\varepsilon$ is a parameter that can be chosen freely. Every choice of $\varepsilon$ gives a value $k$ for the $\varepsilon$-Cluster Similarity, which in turn gives a running time bound. While in Fig. 1 we plotted the best bound that can be achieved with any $\varepsilon$, Fig. 2 shows for each $\varepsilon$ the bound that can be obtained with this $\varepsilon$. It is theoretically clear that the term $4N/\varepsilon$ (yellow) is decreasing in $\varepsilon$ and $2kC$ (red) is increasing (since $k$ increases with $\varepsilon$). The bound (blue) is the sum of these two terms, and we observe that very small and very large choices of $\varepsilon$ often give less good bounds. However, often the blue curve shows some range in which it is rather flat, indicating that there is a large range of $\varepsilon$ that gives comparable bounds. In particular, it seems that the range $0.15 \leq \varepsilon \leq 0.25$ often gives reasonable bounds.

In Fig. 3 we compare the $\varepsilon$-Cluster Similarity ($k$ normalized as a percent of $N$) with another diversity measure: the mean of the covariance matrix of population error. If both measures of diversity were highly correlated, we would expect that for any fixed $\varepsilon$ (points of the same color), points with larger $x$-value would consistently have smaller $y$-value (since $k$ is an inverse measure of diversity). However, in many plots this correlation is spurious at best. It even appears to have opposite signs in some cases. We conclude that $\varepsilon$-Cluster Similarity measures a different aspect of diversity than the mean of the covariance matrix. We also note that, interestingly, the two problems for which there *does* appear

to be a relation between the two measures (compare-string-lengths and mirror-image) are the two problems with boolean error vectors.

## 6    Conclusions

We have introduced and investigated a new measure of population diversity, the $\varepsilon$-Cluster Similarity. We have theoretically proven that large population diversity makes lexicase selection fast, and empirically confirmed that populations are diverse in real-world examples of genetic programming. Thus we have concluded that diverse populations can explain a substantial part of the discrepancy between the general worst-case bound for lexicase selection, and the fast running time in practice.

Naturally, the question arises whether populations in other areas than genetic programming are also diverse with respect to this measure. Moreover, what other consequences does a diverse population have? For example, does it lead to good crossover results? Does it help against getting trapped in local optima? While the intuitive answer to these questions seems to be Yes, it is not easy to pinpoint such statements with rigorous experimental or theoretical results. We hope that our new notion of population diversity can be a means to better understand such questions.

## References

1. Aenugu, S., Spector, L.: Lexicase selection in learning classifier systems. In: Proceedings of the Genetic and Evolutionary Computation Conference. pp. 356–364 (2019)
2. Doerr, B., Johannsen, D., Winzen, C.: Multiplicative drift analysis. Algorithmica **64**(4), 673–697 (2012)
3. Dolson, E., Ofria, C.: Ecological theory provides insights about evolutionary computation. In: Proceedings of the Genetic and Evolutionary Computation Conference Companion. p. 105–106. GECCO '18, Association for Computing Machinery, New York, NY, USA (2018). https://doi.org/10.1145/3205651.3205780, https://doi.org/10.1145/3205651.3205780
4. Helmuth, T., McPhee, N.F., Spector, L.: Effects of Lexicase and Tournament Selection on Diversity Recovery and Maintenance. In: Proceedings of the 2016 on Genetic and Evolutionary Computation Conference Companion. pp. 983–990. ACM (2016), http://dl.acm.org/citation.cfm?id=2931657
5. Helmuth, T., McPhee, N.F., Spector, L.: The impact of hyperselection on lexicase selection. In: Proceedings of the 2016 on Genetic and Evolutionary Computation Conference. pp. 717–724. ACM (2016), http://dl.acm.org/citation.cfm?id=2908851

6. Helmuth, T., McPhee, N.F., Spector, L.: Program synthesis using uniform muta-
   tion by addition and deletion. In: Proceedings of the Genetic and Evolutionary
   Computation Conference. pp. 1127–1134. GECCO '18, ACM, Kyoto, Japan (15-19
   Jul 2018). https://doi.org/10.1145/3205455.3205603, http://doi.acm.org/10.1145/
   3205455.3205603
7. Helmuth, T., Pantridge, E., Spector, L.: On the importance of special-
   ists for lexicase selection. Genetic Programming and Evolvable Machines pp.
   349–373 (2020). https://doi.org/doi:10.1007/s10710-020-09377-2, https://link.
   springer.com/article/10.1007%2Fs10710-020-09377-2
8. Helmuth, T., Spector, L.: General program synthesis benchmark suite. In: GECCO
   '15: Proceedings of the 2015 conference on Genetic and Evolutionary Computation
   Conference. pp. 1039–1046. ACM, Madrid, Spain (11-15 Jul 2015). https://doi.
   org/doi:10.1145/2739480.2754769, http://doi.acm.org/10.1145/2739480.2754769
9. Helmuth, T., Spector, L.: Explaining and exploiting the advantages of down-
   sampled lexicase selection. In: Artificial Life Conference Proceedings. pp. 341–
   349. MIT Press (13-18 Jul 2020). https://doi.org/10.1162/isal_a_00334, https:
   //www.mitpressjournals.org/doi/abs/10.1162/isal_a_00334
10. Helmuth, T., Spector, L., Matheson, J.: Solving uncompromising problems with
    lexicase selection. IEEE Transactions on Evolutionary Computation **19**(5), 630–
    643 (Oct 2015). https://doi.org/10.1109/TEVC.2014.2362729
11. Jansen, T., Zarges, C.: Theoretical analysis of lexicase selection in multi-objective
    optimization. In: Auger, A., Fonseca, C.M., Lourenço, N., Machado, P., Paquete,
    L., Whitley, D. (eds.) Parallel Problem Solving from Nature – PPSN XV. pp.
    153–164. Springer International Publishing, Cham (2018)
12. La Cava, W., Helmuth, T., Spector, L., Moore, J.H.: A probabilistic and multi-
    objective analysis of lexicase selection and epsilon-lexicase selection. Evolution-
    ary Computation **27**(3), 377–402 (2019). https://doi.org/10.1162/evco_a_00224,
    https://arxiv.org/pdf/1709.05394
13. La Cava, W., Orzechowski, P., Burlacu, B., de Franca, F., Virgolin, M., Jin, Y.,
    Kommenda, M., Moore, J.: Contemporary Symbolic Regression Methods and their
    Relative Performance. Proceedings of the Neural Information Processing Systems
    Track on Datasets and Benchmarks **1** (Dec 2021)
14. La Cava, W., Spector, L., Danai, K.: Epsilon-Lexicase Selection for Regression.
    In: Proceedings of the Genetic and Evolutionary Computation Conference 2016.
    pp. 741–748. GECCO '16, ACM, New York, NY, USA (2016). https://doi.org/10.
    1145/2908812.2908898, http://doi.acm.org/10.1145/2908812.2908898
15. Lengler, J.: Drift analysis. In: Theory of Evolutionary Computation, pp. 89–131.
    Springer (2020)
16. Liskowski, P., Krawiec, K., Helmuth, T., Spector, L.: Comparison of Semantic-
    aware Selection Methods in Genetic Programming. In: Proceedings of the Com-
    panion Publication of the 2015 Annual Conference on Genetic and Evolution-
    ary Computation. pp. 1301–1307. GECCO Companion '15, ACM, New York,
    NY, USA (2015). https://doi.org/10.1145/2739482.2768505, http://doi.acm.org/
    10.1145/2739482.2768505
17. Moore, J.M., Stanton, A.: Tiebreaks and diversity: Isolating effects in lexicase
    selection. The 2018 Conference on Artificial Life pp. 590–597 (2018). https://doi.
    org/10.1162/isal_a_00109
18. Orzechowski, P., La Cava, W., Moore, J.H.: Where are we now? A large bench-
    mark study of recent symbolic regression methods. In: Proceedings of the 2018 Ge-
    netic and Evolutionary Computation Conference. GECCO '18 (Apr 2018). https://

doi.org/10.1145/3205455.3205539, http://arxiv.org/abs/1804.09331, tex.ids: orzechowskiWhereAreWe2018a arXiv: 1804.09331

19. Spector, L.: Assessment of problem modality by differential performance of lexicase selection in genetic programming: a preliminary report. In: Proceedings of the fourteenth international conference on Genetic and evolutionary computation conference companion. pp. 401–408 (2012), http://dl.acm.org/citation.cfm?id=2330846

20. Spector, L., Klein, J., Keijzer, M.: The Push3 execution stack and the evolution of control. In: GECCO 2005: Proceedings of the 2005 conference on Genetic and evolutionary computation. vol. 2, pp. 1689–1696. ACM Press, Washington DC, USA (25-29 Jun 2005). https://doi.org/doi:10.1145/1068009.1068292

21. Spector, L., Robinson, A.: Genetic programming and autoconstructive evolution with the push programming language. Genetic Programming and Evolvable Machines **3**(1), 7–40 (Mar 2002). https://doi.org/doi:10.1023/A:1014538503543, http://hampshire.edu/lspector/pubs/push-gpem-final.pdf

22. Vanneschi, L., Castelli, M., Silva, S.: A survey of semantic methods in genetic programming. Genetic Programming and Evolvable Machines **15**(2), 195–214 (2014)