


Unlocking the Secrets of Software Configuration Landscapes—Ruggedness, Accessibility, Escapability, and Transferability *

Mingyu Huang^{†1}, Peili Mao^{†1} and Ke Li²

¹School of Computer Science and Engineering, UESTC, Chengdu, PR China

²Department of Computer Science, University of Exeter, EX4 4QF, Exeter, UK

 k.li@exeter.ac.uk

Abstract: Modern software systems are often highly configurable to tailor varied requirements from diverse stakeholders. Understanding the mapping between configurations and the desired performance attributes plays a fundamental role in advancing the controllability and tuning of the underlying system, yet has long been a dark hole of knowledge due to their black-box nature and the enormous combinatorial configuration space. In this paper, using 86M evaluated configurations from three real-world systems on 32 running workloads, we conducted one of its kind fitness landscape analysis (FLA) for configurable software systems. With comprehensive FLA methods, we for the first time show that: *i*) the software configuration landscapes are fairly rugged, with numerous scattered local optima; *ii*) nevertheless, the top local optima are highly accessible, featuring significantly larger basins of attraction; *iii*) most inferior local optima are escapable with simple perturbations; *iv*) landscapes of the same system with different workloads share structural similarities, which can be exploited to expedite heuristic search. Our results also provide valuable insights on the design of tailored meta-heuristics for configuration tuning; our FLA framework along with the collected data, build solid foundation for future research in this direction.

Keywords: Configurable software systems, software configuration landscape, fitness landscape analysis, local optima network, exploratory data mining.

1 Introduction

Modern software systems have become increasingly sophisticated and highly configurable. For example, the LINUX kernel has 15,000+ options, most of which can take three values. Without considering constraints, this can yield $\approx 3^{15,000}$ configurations—far more than the number of atoms in the universe ($\approx 10^{80}$). On one hand, by tuning various configuration options, stakeholders are empowered with the flexibility to tailor a system according to their requirements on different non-functional performance attributes (e.g., execution time, throughput) [1–3]; yet on the other hand, stakeholders are usually overwhelmed by the complexity of the system under tune and ending up with innate settings or relying on expert consultations. A thorough understanding of the intricate relationship between configurations and performance is essential for effective control of the system under tune, facilitating various downstream tasks such as analysis [3], debugging [4], adaptation [5], optimization [6, 7], and autotuning [8]. However, this is far from trivial in practice due to the black-box nature of software systems.

The fitness landscape metaphor, pioneered by Wright in 1932 [9], is a fundamental concept in evolutionary biology [10], and has been adapted for analyzing input-output responses of black-box systems across various disciplines, e.g., computer science [11], social science [12], chemistry [13], and engineering [14]. It can be envisioned as a (hyper-)surface as formed by *fitness* (i.e., performance

*This manuscript is currently submitted for possible publication. Reviewers can feel free to use this in peer review.

[†]Mingyu Huang and Peili Mao contributed equally to this work.

attribute in our context) values across a high-dimensional configuration space. Each spatial location within this landscape represents a configuration, with its *height* indicating the fitness and optimal configurations sitting on the peaks.

Since a configuration optimization process is akin to navigating uphill towards these peaks, the topography of fitness landscapes (e.g., the location of fitness peaks and their interconnectivity) plays a pivotal role for understanding both system properties and algorithmic behavior. While a plethora of fitness landscape analysis (FLA) [11] methods have been developed over the past decades for exploring diverse landscape properties, little has been known about the topography of software configuration landscapes. This gap can be partially attributed to the enormous configuration space and the expensive cost of performance evaluation, which hinder the construction of high-fidelity landscapes. Additionally, the variety of use cases, a.k.a., workloads, in real-world engineering poses further challenges in drawing general conclusions.

In this paper, by mapping 32 large configuration landscapes of three representative real-world systems, and by leveraging diverse FLA methods, we substantially advance the understanding of the intricacies of black-box configurable software systems from the following four fundamental aspects.

Ruggedness. A fitness landscape can be either unimodal or multimodal, a.k.a., *rugged*, where ruggedness can pose significant challenges for optimizers in locating the global optimum [6]. Configuration landscapes of various software systems are reported to be highly rugged, characterized by local optima [15]. However, this phenomenon is typically observed during the optimization process, which explores only a small fraction of the entire configuration space. In addition, important properties of local optima, such as their spatial distribution and fitness, have remained largely underexplored. This then leads us to our first research question (RQ): **RQ1:** *‘How rugged are the configuration landscapes, and how are local optima distributed across these landscapes?’*. By applying several classic FLA metrics, we demonstrate that local optima are widespread and randomly dispersed in configuration landscapes, contributing to a highly rugged topography (Section 4.1).

Accessibility of global peaks. In such rugged landscapes, featured in tens of thousands of local optima, locating the most prominent peaks might seem as hard as finding a needle in a haystack. While it has been widely known that the probability of encountering a local optimum during optimization depends on its size of basin of attraction, the following fundamental question remains unclear: **RQ2:** *‘whether global peaks in configuration landscapes have tiny, needle-like basins, or are their basins larger than others?’*. Our analyses on over 1M local optima and large-scale simulation, reveals that the ruggedness does not preclude the top local optima to be highly accessible: they typically feature significantly larger basins, which allow them to be reached with a higher probability (Section 4.2).

Escaping inferior peaks. Despite this, it is not uncommon in practice that optimizers can sometimes get trapped by less-fit local optima. Then, a natural follow-up question is *‘whether we can escape such undesired local optima, and thus transition to superior ones’*, which constitutes our **RQ3**. We investigated this by analyzing the local optima network (LON) [16] of each landscape, which is a well-established FLA tool for modeling the interconnectivity pattern between local optima. We observe that inferior local optima can be easily escaped by applying some simple perturbations, and possibly transition to the top peaks (Section 4.3).

Landscape similarity. It is recently shown that landscapes of the same family of black-box optimization problems may share inherent *structural similarities* [17], which provides fundamental support for the effectiveness of transfer learning methods [18]. Here, we hypothesize that this is also true for configurable software systems, and thus our **RQ4:** *‘Do configuration landscapes of the same system with different workloads share any structural similarities?’* By quantifying the consistency of configuration ranks across distinct workloads, and the similarity of learned embeddings for LONs, we find inspiring evidence to support our hypothesis. We further show that such similarities can be exploited to expedite the optimization process (Section 4.4).

To the best of our knowledge, this work represents the first attempt to conduct exploratory analysis on the landscapes of configurable software systems. Towards this end, ► we develop a highly-scalable FLA framework that can manipulate landscapes with millions of configurations (Section 3). ► We

collect extensive performance data from 3 prevalent software systems, with 32 workloads, and over 86M configurations in total (Section 3). ► We for the first time provide insights into four fundamental aspects of configuration landscapes.

2 Background

2.1 Software Configuration Problem

A configurable software system often comes with a set of configuration options $\{c_1, \dots, c_n\}$, each of which can take either categorical (or boolean) or integer values. The whole configuration space, $\mathcal{C} = C_1 \times \dots \times C_n$, is the Cartesian product of the domains of all n options of interest. Let $f : \mathcal{C} \rightarrow \mathbb{R}$ be a fitness function that maps from a configuration vector $\mathbf{c} = (c_1, \dots, c_n) \in \mathcal{C}$ to a performance attribute such as execution time. The goal of software configuration is to find an optimum configuration $\mathbf{c}^* = \underset{\mathbf{c} \in \mathcal{C}}{\operatorname{argmax}} / \underset{\mathbf{c} \in \mathcal{C}}{\operatorname{argmin}} f(\mathbf{c})$. Due to the lack of analytical form of f , this is a black-box optimization problem that is known to be \mathcal{NP} -hard [19].

2.2 Fitness Landscape

Fitness landscape. A *fitness landscape* can be defined as a triplet $(\mathcal{C}, \mathcal{N}, f)$, where \mathcal{C} is the search space; \mathcal{N} indicates a neighborhood structure that specifies the neighbors of each configuration; and $f : \mathcal{C} \rightarrow \mathbb{R}$ is the fitness function.

Local optimum. A configuration \mathbf{c} is said to be a *local optimum* (denoted as \mathbf{c}^ℓ) iff $f(\mathbf{c}^\ell)$ is better than $f(\mathbf{c})$, $\forall \mathbf{c} \in \mathcal{N}(\mathbf{c}^\ell)$, where $\mathcal{N}(\mathbf{c}^\ell)$ is the neighborhood of \mathbf{c}^ℓ .

Basin of attraction. The *basin of attraction* of a local optimum \mathbf{c}^ℓ , denoted as $\mathcal{B}(\mathbf{c}^\ell)$, is the set of all configurations from which local search converges to it, i.e., $\mathcal{B}(\mathbf{c}^\ell) = \{\mathbf{c} \in \mathcal{C} \mid \text{LocalSearch}(\mathbf{c}) \rightarrow \mathbf{c}^\ell\}$. We define the *size* of a basin to be its cardinality $|\mathcal{B}|$, and its average *radius* to be the expected number of local search steps to arrive at the corresponding local optimum. Note that depending on the local search strategy, e.g., *best-* or *first-improvement* [20], the exact basin size and radius can vary.

Local optima network. LON is rooted in the study of energy landscapes in chemical physics [21] and is a compact graph representation of fitness landscape (denoted as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$). It compresses all configurations in the basin of a local optimum \mathbf{c}_i^ℓ into a single vertex $v_i \in \mathcal{V}$ in LON. Each edge $e^{i,j} \in \mathcal{E}$ indicates potential transitions between the local optimum basins, i.e., whether a configuration in $\mathcal{B}(\mathbf{c}_i^\ell)$ can *escape* to another basin $\mathcal{B}(\mathbf{c}_j^\ell)$ by applying some perturbations. Since LON is able to capture various characteristics regarding LOs and their connectivity patterns, it has become one of the most popular methods for FLA [11].

3 Configuration Landscape Construction

This section first introduces our considered software systems in this study, along with the selected options and performance metrics (Table 1). We then delineate our data collection and landscape construction procedure. Due to the stringent page limit, we left more details in APPENDIX 1 and APPENDIX 2.

3.1 Software Systems and Data Collection

Systems under tune. To ensure practicality of our findings, in this paper we investigated three widely used configurable software systems [2], including the ► LLVM compiler, the ► APACHE web server, and the ► SQLITE database. Additionally, to account for varied real-world engineering use cases, we considered different workloads (\mathbf{w}) for each system, which results in a total of 32 different scenarios (see APPENDIX 1).

Table 1. Meta-information of our experiments.

System	Workloads	Options	Configs.	Total Eva.
LLVM	12	20	1.05M	12.58M
APACHE	10	18	1.77M	17.69M
SQLITE	10	16	5.67M	56.65M
TOTAL:	32			86.92M

Feature selection. While each of these systems could have numerous configurable options, not all of them significantly impact performance. Thus, we employed a two-stage feature selection process to identify a more relevant subset of options: ► Initially, we conducted a pre-selection based on each system’s official documentation, including only those options explicitly stated to affect performance. ► Then, we performed an ablation analysis [22] to further filter out options without a statistically significant impact on performance. Following this approach, we finally identified 20 options for LLVM, 18 for APACHE, and 16 for SQLITE.

Performance metrics. In this study, we considered the performance metrics as in the official benchmarking toolkit of each system (see links below). Specifically, we analyze the ► execution time of compiled projects for LLVM[Ⓢ], the ► number of requests handled per second for APACHE[Ⓢ], and the ► number of written items per second for SQLITE[Ⓢ].

Data collection. To collect benchmark data that is compatible with FLA, we exhaustively explored all combinations of categorical options and used a grid search for integer options¹. This approach resulted in 1.7M configurations for each workload of LLVM, 2.9M for APACHE, and 5.0M for SQLITE. For each configuration, we conducted five runs to evaluate the targeted performance metrics and recorded the median value as the final performance indicator. It is important to note that configuration options not considered in this study were set to their recommended values during this data collection phase.

3.2 Construction Method

Distance measures and neighborhood structure. Central to our configuration landscape is a proper notion of *distance* $d(\mathbf{c}_i, \mathbf{c}_j)$ between two configurations where $\mathbf{c}_i, \mathbf{c}_j \in \mathcal{C}$. For categorical (boolean) variables, we employ the Hamming distance as the measure, while for integer variables, the distance is defined as the number of steps between values on their respective grids. Consequently, $d(\mathbf{c}_i, \mathbf{c}_j)$ can be calculated as the sum of the distances for each corresponding pair of values. Given this, the neighborhood $\mathcal{N}(\mathbf{c})$ of a configuration \mathbf{c} is then defined as the set of all configurations differing from \mathbf{c} by exactly one distance, i.e., $\mathcal{N}(\mathbf{c}) = \{\mathbf{c}' \mid d(\mathbf{c}', \mathbf{c}) = 1\}$.

Configuration landscape as a graph. Given the inherent neighborhood structure in fitness landscapes, representing the data as a *graph* is a natural choice. In this graph model, each configuration $\mathbf{c} \in \mathcal{C}$ is represented as a vertex, with the performance metric $f(\mathbf{c})$ assigned as the node attribute. Neighboring configurations $\mathbf{c}' \in \mathcal{N}(\mathbf{c})$, are connected to \mathbf{c} via *directed edge*. The direction of each edge is determined by the relative values of $f(\mathbf{c})$ and $f(\mathbf{c}')$, always pointing towards the *fitter* configuration. This graph-based approach to modeling the configuration landscape then enables the identification of local optima and their basins, as well as the implementation of FLA methods, using straightforward graph traversal techniques (see APPENDIX 2). This allows our FLA framework to efficiently explore probably the largest landscapes ever in existing literature, comprising millions of configurations.

LON construction. The essence of LON construction lies in defining the transitions between local optima, following the established routine in [23] (Algorithm 1). Specifically, for each local optimum $\mathbf{c}^\ell \in \mathcal{C}$, we apply a k -kick ($k \geq 2$) perturbation to it, resulting in a new configuration \mathbf{c}'' where $d(\mathbf{c}^\ell, \mathbf{c}'') = k$. A local search initiated from \mathbf{c}'' leads to a new local optimum $\mathbf{c}_{\text{new}}^\ell$. If $\mathbf{c}_{\text{new}}^\ell$ is identical to \mathbf{c}^ℓ , the process fails to *escape* from \mathbf{c}^ℓ . Otherwise, an edge is drawn from \mathbf{c}^ℓ to $\mathbf{c}_{\text{new}}^\ell$. In this study,

¹All relevant data can be found in <http://tinyurl.com/4vukrtct>.

Algorithm 1: Constructing local optima network

Input: The set of local optima \mathcal{V} ; Configuration space \mathcal{C} ; A neighborhood function $\mathcal{N}_d(\mathbf{c})$

Output: $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$

```
1 foreach  $\mathbf{c}^\ell \in \mathcal{V}$  do
2   foreach  $\mathbf{c}'' \in \mathcal{N}_2(\mathbf{c}^\ell)$  do
3      $\mathbf{c}_{\text{new}}^\ell \leftarrow \text{LOCALSEARCH}(\mathbf{c}'')$ ;
4     if  $f(\mathbf{c}_{\text{new}}^\ell) < f(\mathbf{c}^\ell)$  and  $\mathbf{c}_{\text{new}}^\ell \neq \mathbf{c}^\ell$  then
5        $\mathcal{E} \leftarrow \mathcal{E} \cup \{(\mathbf{c}^\ell, \mathbf{c}_{\text{new}}^\ell)\}$ ;  $\mathcal{W}[(\mathbf{c}^\ell, \mathbf{c}_{\text{new}}^\ell)] \leftarrow 1$ ;
6     else
7        $\mathcal{W}[(\mathbf{c}^\ell, \mathbf{c}_{\text{new}}^\ell)] \leftarrow \mathcal{W}[(\mathbf{c}^\ell, \mathbf{c}_{\text{new}}^\ell)] + 1$ ;
```

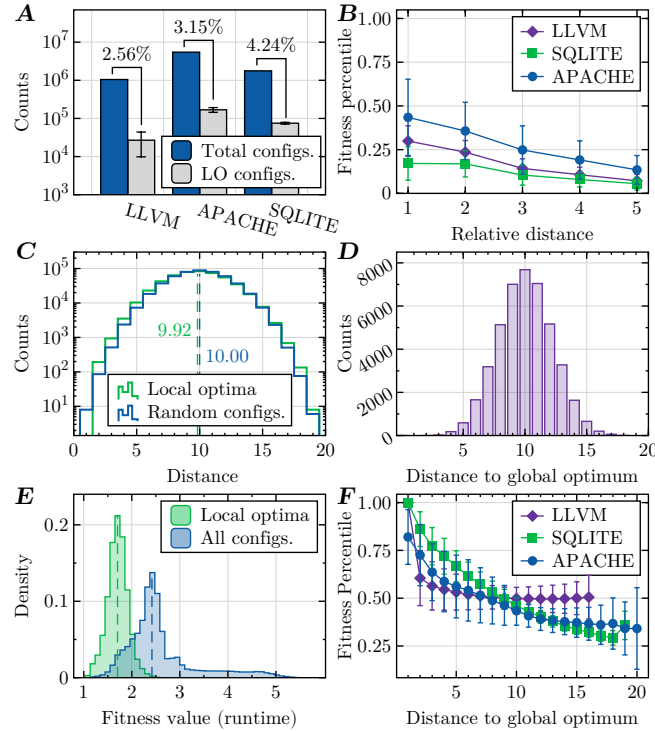


Figure 1. (A) Number of total configurations and local optima in each system, aggregated across workloads. (B) Autocorrelation calculated under different distances, aggregated across workloads. (C) Distribution of pairwise distance between local optima and 10^4 randomly sampled configurations for LLVM-w1, where dashed lines are means. (D) Distribution of the distance of each local optimum to the global optimum for LLVM-w1. (E) Distribution of fitness values of all configurations local optima for LLVM-w1. (F) Average local optimum fitness (in percentile, lower is better) versus the corresponding distance to the global optimum for each system.

we examine all 2-kick perturbations for each local optimum and record transition frequencies as edge weights. Additionally, following [24], we retain only *improving edges*, i.e., those leading to a local optimum with better fitness, which results in a so-called *monotonic* LON.

4 Configuration Landscape Analysis

This section addresses the RQs outlined in Section 1 with a series of dedicated FLA methods. For sake of brevity, some per-instance discussions will focus on the LLVM-w1. Note that this is not indicative of bias: as detailed in APPENDIX 3, our findings are representative across all 32 scenarios.

4.1 Configuration Landscapes are Highly Rugged

In order to investigate **RQ1**, we evaluated the number of local optima in each landscape and the autocorrelation metric. We also analyzed the spatial distribution of local optima as well as their fitness values, and thereby offer a comprehensive assessment of the landscape’s ruggedness.

Number of local optima. We first quantified the number of local optima in each constructed landscape, which serves as a coarse-grained indicator of landscape ruggedness. We found that all investigated landscapes exhibit a significant number of local optima, ranging from 10^4 to 10^5 , which take on average 2.56% to 4.24% of the total configurations (see Fig. 1A). This indicates a high degree of ruggedness in the software configuration landscapes that is comparable to the maximally rugged Kauffman’s NK -landscapes [25].

Autocorrelation. To further validate the observed ruggedness, we calculated the *autocorrelation* metric [26] for each landscape, which is a widely used method for rigorously assessing landscape ruggedness. This metric evaluates the extent to which configurations in close proximity (in terms of distance) tend to exhibit similar fitness values. It is defined by the autocorrelation of fitness values across configurations visited during random walks through the landscape. The corresponding results show a high concordance with the large number of local optima (see Fig. 1B), indicating that even in the most local regions (i.e., $d = 1$), the fitness of configurations is only weakly correlated. This correlation diminishes further as the distance between configurations increases.

Distribution of local optima. In landscapes with numerous local optima, understanding their distribution is crucial. In particular, we seek to ascertain whether they are closely clustered or widely dispersed. To investigate this, we analyzed the distribution of pairwise distances between local optima and compared it with that of 10K randomly sampled configurations from each landscape. Surprisingly, across all scenarios, the two distributions were remarkably similar (two-sided Kolmogorov-Smirnov test $D < 0.147$, $p > 0.954$; see Fig. 1C). This suggests that local optima are almost uniformly distributed across the landscapes. While this analysis considers distances between all pairs of local optima, we also examined particular regions of interest. For example, Fig. 1D depicts the distribution of distances between each local optimum and the global optimum in LLVM-W1, which reveals that most local optima are located far from the global peak, following a normal distribution with a median distance of 10—the same as the radius of the LLVM landscape.

Fitness distance correlation. In addition to their spatial locations, the fitness distribution of local optima is also of significant importance in rugged landscapes. For instance, landscapes where the fitness of a local optimum negatively correlates with its distance to the global optimum can be ‘deceptive’ [27], potentially misleading the optimizer. In our cases, we found that most local optima in configuration landscapes are fitter than random configurations, but only a few approach the fitness of the optimum (see Fig. 1E). In general, sub-optimal local optima are located in the vicinity of the global optimum, and their fitness generally decreases with increasing distance (see Fig. 1F). However, such correlation, a.k.a., fitness distance correlation [28], is weak for most cases ($|r| < 0.2$). Thus, it is not unusual in these landscapes to find moves that increase fitness but paradoxically lead further away from the true optimum.

💡 **Response to RQ1:** *Software configuration landscapes are highly rugged, where local optima are prevalent and nearly randomly dispersed. Most of them are far away from the global optimum, and as they approach the optimum, their fitness tend to increase.*

4.2 Global Peaks are Highly Accessible

Addressing **RQ2** is not straightforward, since there exist two intertwined paradigms for identifying basins of a local optimum: the *first*- and *best*-improvement local search. We conducted analyses under both paradigms to see whether top local optima possess larger basins, thereby increasing their likelihood of being encountered during the search process.

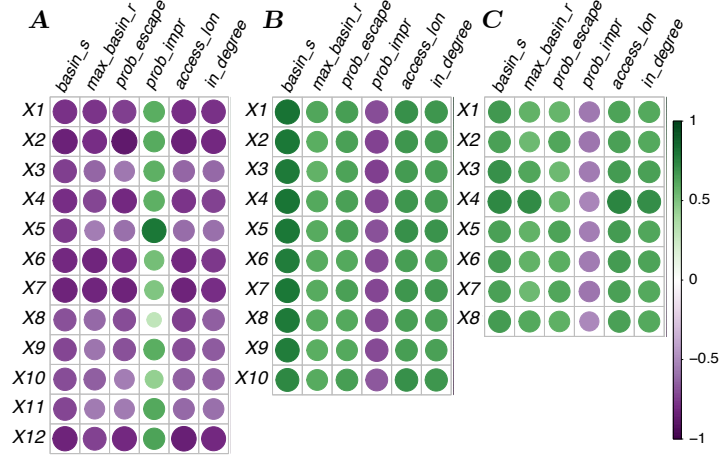


Figure 2. Spearman correlation between local optima fitness and their properties for (A) LLVM, (B) SQLite, and (C) APACHE.

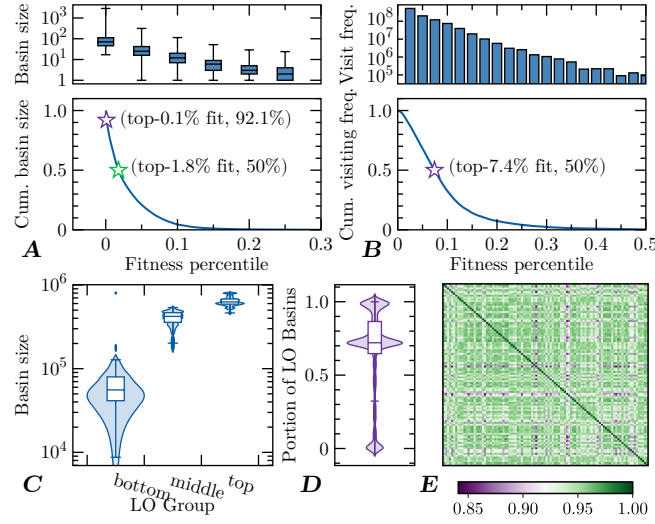


Figure 3. Example results on LLVM-w1: (A) boxplots and cumulative distribution of the size of best-improvement basins versus fitness percentile (the lower the better). (B) bars and cumulative distribution of frequency of visits for local optima at different fitness percentiles, during 10^9 first-improvement local search runs. The long tails of both curves have been truncated. (C) distribution of basin size for local optima with fitness in the top, middle and bottom 0.1% percentile based on first-improvement local search. (D) the distribution of accessibility to local optima basins of 10^3 random configurations, where 1.0 implies a configuration is shared by all basins. (E) overlap ratio between basins of the top 100 local optima.

Best-improvement basins. In order to determine the basins of best-improvement for each local optimum, we conducted exhaustive best-improvement local searches starting from every configuration within the landscapes, and recorded the terminating local optimum. For all 32 landscapes investigated, we observed significant correlation between basin size and the fitness of local optima (see Fig. 2). In particular, local optima with lower fitness have smaller basins, encompassing only a few configurations. On the other hand, those with higher fitness tend to feature significantly larger basins, ranging from hundreds to thousands of configurations (see Fig. 3A). Considering that under best-improvement search, each configuration deterministically converges to a specific basin, the expected basin size for a local optimum in our landscapes is thereby relatively small, often just a few dozen of configurations. This highlights that global peaks’ basins are substantially larger than average, increasing their accessibility during optimization processes. Actually, since the basin of each local optimum is exhaustively determined, the likelihood of reaching a given local optimum via a randomly initiated best-improvement search is then directly proportional to its basin size. By plotting the cumulative distribution of local optimum basin sizes, plotted against fitness percentiles, we found that there is a

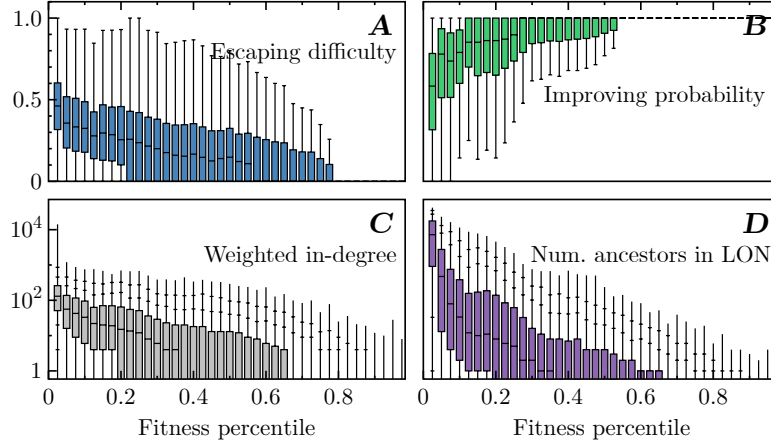


Figure 4. Distribution of local optimum attributes calculated from LONs versus fitness percentile (the lower the better) for LLVM-W1.

50% probability of achieving a satisfactory fitness level at the first encounter of a local optimum (see ☆ in Fig. 3A; APPENDIX 3). Even more notably, for the top 0.1% of local optima, there remains a competitive probability of reaching them on the first attempt (see ☆ in Fig. 3A).

First-improvement basins. Our previous observations also persist if we consider the first-improvement paradigm. Here, despite significant speed-up strategies in our FLA framework, it is still computationally prohibitive to exhaustively identify the first-improvement basins. Alternatively, we focused on analyzing basins for local optima within the top, middle, and bottom 0.1% fitness percentiles, respectively. The results revealed that the top local optima possess dominantly large basins, with their median encompassing over half of all configurations, significantly surpassing the other two groups (see Fig. 3C, Wilcoxon rank-sum test $p < 10^{-i}$). This observation led us to hypothesize substantial overlap among the largest basins. Indeed, we found that the basins of the top 0.1% local optima overlap by an average of $94.4\% \pm 2.1\%$ (mean \pm standard deviation). Moreover, an inspection of 10^3 random configurations in each landscape showed that most simultaneously lie within the basins of more than 50% of the local optima (see Fig. 3D). Such overlap can then give rise to optimization contingency, where the same starting configuration might lead to different local optima, not necessarily be the best accessible one. To investigate this, and to verify the accessibility of global peaks, we conducted 10^9 runs of first-improvement local search across each landscape. We observed that fitter local optima were visited significantly more often (see Fig. 3B; APPENDIX 3), with half of the runs reaching local optima within the top 7.4% fitness percentile (☆). While notable, these results are less pronounced compared to those from greedy best-improvement search, likely due to the extensive basin overlap [20].

💡 **Response to RQ2:** *Global peaks in software configuration landscapes feature dominant size of basins, which allow them to be accessed with higher chance.*

4.3 Local Optima Can be Easily Escaped

To address **RQ3**, we need an appropriate model to model the potential connectivity among different local optimum basins. Accordingly, we constructed an LON for each landscape as outlined in Section 3. Subsequently, we analyzed various features of local optima extracted from these LONs.

Escaping from inferior local optima. To order to explore how challenging it is to escape from a local optimum in configuration landscapes, we measured the *escaping difficulty* for each local optimum within the LON. This metric is defined by the proportion of 2-kick perturbations that return to the original local optimum. As shown in Fig. 2, there is a negative correlation between escaping difficulty and the fitness of local optima. Specifically, the least fit local optima can be escaped from using almost any 2-kick perturbation (see Fig. 4A). We attribute this to their smaller basin sizes, a hypothesis supported by a strong correlation between basin size and escaping difficulty (0.90 ± 0.04). As local

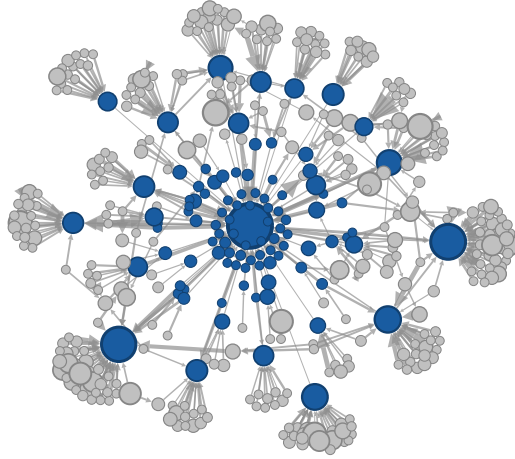


Figure 5. A part of the LLVM-W1 LON containing the global peaks. The global optimum and its nearest neighbors are colored in blue, and node radius indicates the basin size.

optima become fitter, escaping difficulty increases, yet even near-optimal ones offer a considerable chance of escape, which benefits the exploration phase of optimization algorithms. However, escaping from a local optimum does not always lead to improved fitness. To validate this assertion, we calculated the *improving probability* for each local optimum, defined as the fraction of escape moves leading to a fitter local optimum. Interestingly, this probability also shows a negative correlation with fitness values (Fig. 2). While nearly all escape moves from inferior local optima lead to higher fitness, this likelihood decreases for fitter local optima (Fig. 4B). Nevertheless, even for top local optima, about half of the escape moves could potentially result in improvements.

Transitions to global peaks. A pertinent question arising from our analysis is whether such escapability can eventually lead to transitions from local optima to global peaks. To find out, we first examined the *in-degree* of each local optimum within the LON, weighted by the transition frequency. The results indicate that top local optima generally possess higher in-degrees than their inferior counterparts (see Fig. 4C), suggesting a relative ease in transitioning from less optimal regions to more advantageous ones, even with simple 2-kick perturbations. While this demonstrates direct connections between local optima, it does not capture the entire connectivity landscape. Therefore, we also analyzed the total number of *ancestors* for each local optimum in the LON. This analysis revealed that, on average, $78.4\% \pm 16.5\%$ of local optima can transition to the top-10 local optima across all scenarios. Fig. 5 depicts the topology of one such LON, illustrating that despite the geographical randomness of local optima, there is a recognized pattern of interconnectedness forming local communities, or landscape *funnels* [29]. These funnels are typically centered around prominent local optima with large basins, surrounded by numerous less fit local optima capable of transitioning towards them.

To verify these findings, we ran two classic search-based software engineering [30] algorithms, particle swarm optimization (PSO) and genetic algorithm (GA), on all the landscapes, with a budget of 2,000 function evaluations and 100 repetitions. The results reported in Fig. 5 shows that both algorithms can consistently locate the top 0.02% local optima, with some even approaching the 0.001% level, nearly approaching the global optimum.

💡 **Response to RQ3:** *Despite the ruggedness of software configuration landscapes, inferior local optima can be easily escaped with simple perturbations, and can largely transition to the global peaks.*

4.4 Landscape Similarities Across Workloads

To investigate RQ4, we first assessed the similarity in configuration rankings across different workloads within the same system. We further adopted high-level features learned from LONs as proxies to

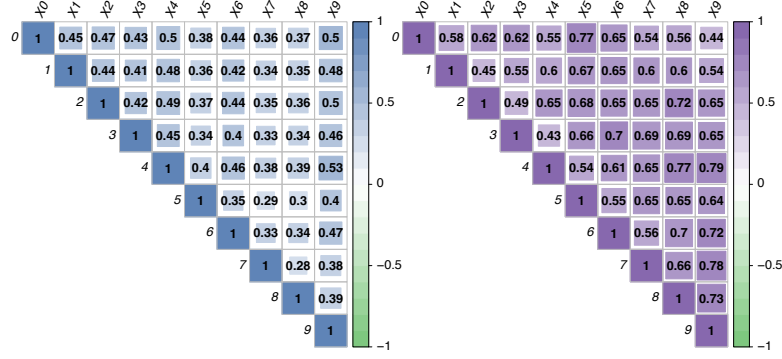


Figure 6. Similarity between the configuration landscapes of LLVM across workloads as measured by (left): Spearman correlation between `run_time` of configurations, and (right): structural similarity of the corresponding LONs.

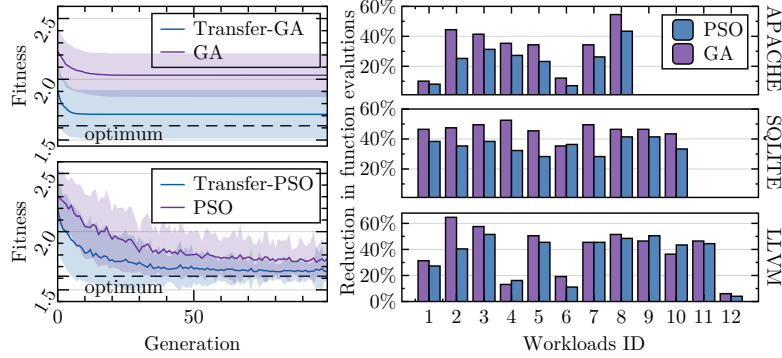


Figure 7. (left): evolution trajectories of warm-started and vanilla version of GA and PSO on LLVM-w1. (right): reduction in function evaluations for warm-started version of GA and PSO to reach the best fitness level achieved by the vanilla version on LLVM-w1.

quantify structural similarities between landscapes of different workloads. Finally, we demonstrated how these identified similarities can be effectively utilized to expedite optimization via warm-start [31].

Similarity across workloads. We first leveraged the Spearman correlation coefficient to quantify the similarity between fitness values of different workloads within each system. This served as an indicator of whether a configuration tends to have similar performance ranks for different workloads. The results show that, while coefficient values can vary, in general a positive correlation exists between fitness values across workloads (see Fig. 6 left). In addition, inspired by the recent work [17], we applied the `GL2Vec` graph embedding method [32] to generate high-level feature representations of each LON, which are able to preserve the intrinsic topology of the network. From the Spearman correlation of these vectors shown in Fig. 6 (right), we clearly see that, the LON structure of most workloads are highly similar. This aligns with the fitness correlation, and implies that the connections and clustering of local optima follow similar patterns across various workloads.

Warm-starting optimizations. The observed fitness correlation and structural similarity make it natural and appealing to warm-start an optimization process by reusing configurations that are known to perform well on similar tasks. To validate this assertion, we integrated such mechanism into PSO and GA. Specifically, for each workload, we initialized the population using the best configurations found in historical optimizations of the 3 most similar tasks. In particular, the historical optimization data were readily available from **RQ3**, and we leverage the structural similarities between LONs as the measure of task similarity. We ran all both warm-started and vanilla version of PSO and GA on each workload for 100 times to ensure statistically robust conclusions. As depicted in the left panel of Fig. 7 (full results are in APPENDIX 4), we found that the warm-started version of both algorithms can outperform the vanilla version in most cases. More importantly, by employing warm-start, we typically need considerably fewer function evaluations to reach the best fitness level achieved by the vanilla version—As can be seen from Fig. 7 (right), in most scenarios the warm-started algorithms

can save over 20% function evaluations, and in some cases even up to 60%.

💡 **Response to RQ4:** *There exists structural similarity between configuration landscapes of different workloads, which can be exploited to expedite optimization via reusing prominent configurations in similar tasks for warm-start.*

5 Related Work

Software configuration tuning. Efficiently tuning of configurable software systems has been a long-desired goal in software engineering. Although a plethora of techniques have been developed to automate this laborious process, e.g., random search [33], hill-climbing [34], genetic algorithms [35], and sequential model-based ones [7], they are usually based on intuition or certain assumption of the configuration space. Consequently, they may not be sufficiently efficient due to the lack of knowledge for the underlying system.

Performance analysis on software systems. Traditionally, the configuration-performance interaction has been explored by *black-box* approaches like performance-influence modeling [2, 36, 37], where the goal is to model such interaction using representative samples drawn from the system under investigation. Another line of related research use machine learning methods to predict the performance of a given configuration [7, 38, 39]. While our collected data can also be used for the above purposes, we note that they are however, orthogonal and complementary to this study. For example, [40] found that performance models can be linearly transferred between workloads, which aligns with our results for **RQ4**, but we explore this from a landscape topology perspective. On the flip side, there is a growing interest in applying *white-box* approaches to explore additional insights of the interactions between configuration options versus performance, e.g., [3, 41]. Particularly, the white-box methods focus on interrogating the system’s source codes to identify the regions of a system responsible for the performance difference among configurations, which is also orthogonal to this present work.

Landscape analysis for software system. Landscape analysis has been a subject of rigorous research within the meta-heuristics community for several decades, but is a fairly new topic in software engineering. Some works have briefly discussed the local optima in software configuration [42, 43], but they achieved this only with naive 2/3D projections. So far there is no dedicated work on comprehensively investigating the topography of configuration landscapes, and this work is the first to do so.

6 Conclusion and Discussions

This work conducted the first FLA on 32 configuration landscapes of three real-world software systems. We observed a fairly rugged topography for these landscapes with abundant pervasive local optima (**RQ1**), which can potentially present challenge for not only navigation, but also performance modeling, as have been known in the protein engineering community [44]. Intriguingly, we later observed that albeit their rugged surface, configuration landscapes are highly navigable in the sense that global peaks often feature dominantly large basins of attraction (**RQ2**), and those inferior ones can be easily escaped with simple perturbations (**RQ3**). These then highlight the need for efficient exploration and escaping mechanisms for searching. Finally, we discovered structural similarity among configuration landscapes across workloads (**RQ4**), providing concrete evidence to advocate configuration reuse strategies in practice, which has attracted attention in the community [45].

There are far more interesting properties regarding configuration landscape that have yet to be explored in the future. We hope that this work will be the first of a continually growing sequence of rigorous research for configuration landscape analysis. Additionally, in open-sourcing the developed FLA framework and gathered performance data, we also hope to make landscape analysis and software configuration research more accessible and reproducible by the community.

Acknowledgment

This work was supported in part by the UKRI Future Leaders Fellowship under Grant MR/S017062/1 and MR/X011135/1; in part by NSFC under Grant 62376056 and 62076056; in part by the Royal Society under Grant IES/R2/212077; in part by the EPSRC under Grant 2404317; in part by the Kan Tong Po Fellowship (KTP\R1\231017); and in part by the Amazon Research Award and Alan Turing Fellowship.

References

- [1] T. Xu, L. Jin, X. Fan, Y. Zhou, S. Pasupathy, and R. Talwadker, “Hey, you have given me too many knobs!: understanding and dealing with over-designed configuration in system software,” in *ESEC/FSE’15*. ACM, 2015, pp. 307–319. 1
- [2] N. Siegmund, A. Grebhahn, S. Apel, and C. Kästner, “Performance-influence models for highly configurable systems,” in *ESEC/FSE’15*. ACM, 2015, pp. 284–294. 1, 3, 11
- [3] M. Velez, P. Jamshidi, F. Sattler, N. Siegmund, S. Apel, and C. Kästner, “Configcrusher: towards white-box performance analysis for configurable systems,” *Autom. Softw. Eng.*, vol. 27, no. 3, pp. 265–300, 2020. 1, 11
- [4] X. Han, T. Yu, and D. Lo, “Perflearner: learning from bug reports to understand and generate performance test frames,” in *ASE’18*. ACM, 2018, pp. 17–28. 1
- [5] T. Chen, K. Li, R. Bahsoon, and X. Yao, “FEMOSAA: feature-guided and knee-driven multi-objective optimization for self-adaptive software,” *ACM Trans. Softw. Eng. Methodol.*, vol. 27, no. 2, pp. 5:1–5:50, 2018. 1
- [6] Y. Zhu, J. Liu, M. Guo, Y. Bao, W. Ma, Z. Liu, K. Song, and Y. Yang, “Bestconfig: tapping the performance potential of systems via automatic configuration tuning,” in *SoCC’17*. ACM, 2017, pp. 338–350. 1, 2
- [7] V. Nair, Z. Yu, T. Menzies, N. Siegmund, and S. Apel, “Finding faster configurations using FLASH,” *IEEE Trans. Software Eng.*, vol. 46, no. 7, pp. 794–811, 2020. [Online]. Available: <https://doi.org/10.1109/TSE.2018.2870895> 1, 11
- [8] A. H. Ashouri, W. Killian, J. Cavazos, G. Palermo, and C. Silvano, “A survey on compiler autotuning using machine learning,” *ACM Comput. Surv.*, vol. 51, no. 5, pp. 96:1–96:42, 2019. 1
- [9] S. Wright, “The roles of mutations, inbreeding, crossbreeding and selection in evolution,” in *Proc. of the 11th International Congress of Genetics*, vol. 1, 1932, pp. 356–366. 1
- [10] E. D. Vaishnav, C. G. de Boer, J. Molinet, M. Yassour, L. Fan, X. Adiconis, D. A. Thompson, J. Z. Levin, F. A. Cubillos, and A. Regev, “The evolution, evolvability and engineering of gene regulatory DNA,” *Nature*, vol. 603, no. 7901, pp. 455–463, 2022. [Online]. Available: <https://doi.org/10.1038/s41586-022-04506-6> 1
- [11] K. M. Malan, “A survey of advances in landscape analysis for optimisation,” *Algorithms*, vol. 14, no. 2, p. 40, 2021. 1, 2, 3
- [12] D. A. Levinthal, “Organizational adaptation and environmental selection-interrelated processes of change,” *Org. Sci.*, vol. 2, no. 1, pp. 140–145, 1991. 1
- [13] J. P. Doye, “Network topology of a potential energy landscape: A static scale-free network,” *Phy. Rev. Lett.*, vol. 88, no. 23, p. 238701, 2002. 1

- [14] V. K. Vassilev, J. F. Miller, and T. C. Fogarty, “Digital circuit evolution and fitness landscapes,” in *CEC’99*. IEEE, 1999, pp. 1299–1308. 1
- [15] P. Jamshidi and G. Casale, “An uncertainty-aware approach to optimal configuration of stream processing systems,” in *MASCOTS’16*. IEEE Computer Society, 2016, pp. 39–48. 2
- [16] G. Ochoa, M. Tomassini, S. Vérel, and C. Darabos, “A study of NK landscapes’ basins and local optima networks,” in *GECCO’08*. ACM, 2008, pp. 555–562. 2
- [17] M. Huang and K. Li, “Exploring structural similarity in fitness landscapes via graph data mining: A case study on number partitioning problems,” in *IJCAI’23*. ijcai.org, 2023, pp. 5595–5603. [Online]. Available: <https://doi.org/10.24963/ijcai.2023/621> 2, 10
- [18] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, 2010. 2
- [19] T. Weise, “Global optimization algorithms-theory and application,” *Self-Published Thomas Weise*, vol. 361, 2009. 3
- [20] L. D. Whitley, A. E. Howe, and D. Hains, “Greedy or not? best improving versus first improving stochastic local search for MAXSAT,” in *AAAI’13*. AAAI Press, 2013, pp. 940–946. 3, 8
- [21] F. H. Stillinger, “A topographic view of supercooled liquids and glass formation,” *Science*, vol. 267, no. 5206, pp. 1935–1939, 1995. 3
- [22] A. Biedenkapp, M. Lindauer, K. Eggensperger, F. Hutter, C. Fawcett, and H. H. Hoos, “Efficient parameter importance analysis via ablation with surrogates,” in *AAAI’17*. AAAI Press, 2017, pp. 773–779. 4
- [23] G. Ochoa and N. Veerapen, “Deconstructing the big valley search space hypothesis,” in *EvoCOP’16*, vol. 9595. Springer, 2016, pp. 58–73. 4
- [24] S. L. Thomson, F. Daolio, and G. Ochoa, “Comparing communities of optima with funnels in combinatorial fitness landscapes,” in *GECCO’17*. ACM, 2017, pp. 377–384. 5
- [25] S. A. Kauffman, *The Origins of Order: Self-Organization and Selection in Evolution*. Oxford University Press, 1993. 6
- [26] E. Weinberger, “Correlated and uncorrelated fitness landscapes and how to tell the difference,” *Biol. Cybern.*, vol. 63, no. 5, pp. 325–336, 1990. 6
- [27] M. Mitchell, *An Introduction to Genetic Algorithms*. MIT Press, 1998. 6
- [28] T. Jones and S. Forrest, “Fitness distance correlation as a measure of problem difficulty for genetic algorithms,” in *ICGA’95*. Morgan Kaufmann, 1995, pp. 184–192. 6
- [29] J. N. Onuchic and P. G. Wolynes, “Theory of protein folding,” *Current opinion in structural biology*, vol. 14, no. 1, pp. 70–75, 2004. 9
- [30] M. Harman, S. A. Mansouri, and Y. Zhang, “Search-based software engineering: Trends, techniques and applications,” *ACM Comput. Surv.*, vol. 45, no. 1, pp. 11:1–11:61, 2012. 9
- [31] M. Feurer, J. T. Springenberg, and F. Hutter, “Initializing bayesian hyperparameter optimization via meta-learning,” in *AAAI’15*. AAAI Press, 2015, pp. 1128–1135. 10
- [32] H. Chen and H. Koga, “Gl2vec: Graph embedding enriched by line graphs with edge features,” in *ICONIP’19*, ser. Lecture Notes in Computer Science, T. Gedeon, K. W. Wong, and M. Lee, Eds., vol. 11955. Springer, 2019, pp. 3–14. 10

- [33] J. Oh, D. S. Batory, M. Myers, and N. Siegmund, “Finding near-optimal configurations in product lines by random sampling,” in *ESEC/FSE’17*. ACM, 2017, pp. 61–71. [11](#)
- [34] M. Li, L. Zeng, S. Meng, J. Tan, L. Zhang, A. R. Butt, and N. C. Fuller, “MRONLINE: mapreduce online performance tuning,” in *HPDC’14*. ACM, 2014, pp. 165–176. [11](#)
- [35] A. Shahbazian, S. Karthik, Y. Brun, and N. Medvidovic, “equal: informing early design decisions,” in *ESE/FSE’20*, 2020, pp. 1039–1051. [11](#)
- [36] R. Olaechea, D. Rayside, J. Guo, and K. Czarnecki, “Comparison of exact and approximate multi-objective optimization for software product lines,” in *SPLC’18*. ACM, 2014, pp. 92–101. [11](#)
- [37] C. Kaltenecker, A. Grebhahn, N. Siegmund, J. Guo, and S. Apel, “Distance-based sampling of software configuration spaces,” in *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*. IEEE / ACM, 2019, pp. 1084–1094. [11](#)
- [38] H. Ha and H. Zhang, “Deepperf: performance prediction for configurable software with deep sparse neural network,” in *ICSE’19*. IEEE / ACM, 2019, pp. 1095–1106. [11](#)
- [39] P. Valov, J. Guo, and K. Czarnecki, “Empirical comparison of regression methods for variability-aware performance prediction,” in *SPLC’15*. ACM, 2015, pp. 186–190. [11](#)
- [40] P. Jamshidi, N. Siegmund, M. Velez, C. Kästner, A. Patel, and Y. Agarwal, “Transfer learning for performance modeling of configurable systems: an exploratory analysis,” in *ASE’17*. IEEE Computer Society, 2017, pp. 497–508. [11](#)
- [41] M. Velez, P. Jamshidi, N. Siegmund, S. Apel, and C. Kästner, “White-box analysis over machine learning: Modeling performance of configurable systems,” in *ICSE’21*, 2021, accepted for publication. [11](#)
- [42] P. Jamshidi, M. Velez, C. Kästner, and N. Siegmund, “Learning to sample: exploiting similarities across environments to learn performance models for configurable systems,” in *ESEC/FSE’18*. ACM, 2018, pp. 71–82. [11](#)
- [43] M. J. V. D. Donckt, D. Weyns, F. Quin, J. V. D. Donckt, and S. Michiels, “Applying deep learning to reduce large adaptation spaces of self-adaptive systems with multiple types of goals,” in *SEAMS’20*. ACM, 2020, pp. 20–30. [11](#)
- [44] S. A. Fahlberg, C. R. Freschlin, P. Heinzelman, and P. A. Romero, “Neural network extrapolation to distant regions of the protein fitness landscape,” *bioRxiv*, pp. 2023–11, 2023. [11](#)
- [45] C. Kinneer, D. Garlan, and C. L. Goues, “Information reuse and stochastic search: Managing uncertainty in self-^{*} systems,” *ACM Trans. Auton. Adapt. Syst.*, vol. 15, no. 1, pp. 3:1–3:36, 2021. [11](#)

A Data Collection

A.1 Configurable Software Systems

To ensure the practicality and generality of our empirical findings, this paper considers investigating three widely used configurable software systems with diverse engineering functionalities, including compiler, Web server, and database management system. In the following paragraphs, we outline their key characteristics including the engineering narration, configuration options considered in our experiments, and the settings of workloads.

- **LLVM**²: The LLVM Project is a collection of modular compiler and toolchain technologies. It provides a modern, SSA-based compilation strategy that supports both static and dynamic compilation of any programming language. ► LLVM has more than 578 configuration options while we choose 20 of them for our empirical study. ► 12 test suites from the widely used PolyBench benchmark suite³ are chosen to constitute the workloads. ► The `run_time` for the compiled program is used as the fitness function to measure the quality of a configuration of the LLVM.
- **Apache**⁴: The APACHE HTTP Server Project aims to provide a robust and scalable HTTP service. ► It consists of multiple modules, the core of which has 89 configuration options and 21 configuration options for MPM module. Here we choose 15 options directly related to the quality of a configuration in our experiments. ► 9 different running environments are generated by using the Apache HTTP server benchmarking tool⁵. ► We use the request handled per second as the fitness function.
- **SQLite**⁶: This is an embedded database project. Instead of maintaining a separate server process, SQLITE directly reads and writes data to disk. ► It has 50 compile-time and 29 run-time configuration options and we chose 18 of them in this study. ► We used the SQLITE Benchmark⁷ to constitute 10 different running workloads. ► The writing speed in sequential key order in async mode (`fillseqsync`) is used as the fitness function.

The meta information of the selected configuration options (parameters) for these systems are listed in Tables 2 to 4. The settings of different workloads for each system are listed in Table 5. In the following paragraphs, we introduce the corresponding attributes for different workloads.

- For the LLVM, we adopted different compiling file to constitute different workloads, as indicated by the attribute `program_name`.
- For the APACHE, there are two attributes to setup the system, and their different combinations constitute different workloads:
 - `requests` represents the number of requests to perform for the benchmarking session. The default is to just perform a single request which usually leads to non-representative benchmarking results.
 - `concurrency` represents the number of multiple requests to perform concurrently. The default is one request at a time.
- For the SQLITE, the workloads are based on two system attributes:
 - `num` indicates the number of entries.
 - `value_size` represents the value size.

²<https://llvm.org/>

³<http://web.cse.ohio-state.edu/~pouchet.2/software/polybench/>

⁴<https://httpd.apache.org/>

⁵<https://httpd.apache.org/docs/2.4/programs/ab.html>

⁶<https://www.sqlite.org/index.html>

⁷<https://github.com/ukontainer/sqlite-bench>

Table 2. Selected configuration options for LLVM

Index	Parameter	Value
1	inline	{on, off}
2	openmpopt	{on, off}
3	mldst-motion	{on, off}
4	gvn	{on, off}
5	jump-threading	{on, off}
6	correlated-propagation	{on, off}
7	elim-avail-extern	{on, off}
8	tailcallelim	{on, off}
9	constmerge	{on, off}
10	dse	{on, off}
11	slp-vectorizer	{on, off}
12	callsite-splitting	{on, off}
13	argpromotion	{on, off}
14	aggressive-instcombine	{on, off}
15	polly-simplify	{on, off}
16	polly-dce	{on, off}
17	polly-optree	{on, off}
18	polly-delicm	{on, off}
19	polly-opt-is1	{on, off}
20	polly-prune-unprofitable	{on, off}

Table 3. Selected configuration options for APACHE

Index	Parameter	Value
1	AcceptFilter	{nntp, http}
2	KeepAlive	{on, off}
3	KeepAliveTimeout	{1, ..., 300}
4	MaxKeepAliveRequests	{1, ..., 2 ¹⁰ }
5	TimeOut	{1, ..., 300}
6	MaxConnectionsPerChild	{1, ..., 1,000}
7	MaxMemFree	{2 ¹⁰ , ..., 2 ²⁰ }
8	MaxRequestWorkers	{100, ..., 3,000}
9	MaxSpareThreads	{50, ..., 500}
10	MinSpareThreads	{20, ..., 250}
11	SendBufferSize	{2 ¹⁰ , ..., 2 ¹⁶ }
12	ServerLimit	{100, ..., 3,000}
13	StartServers	{1, ..., 10}
14	ThreadLimit	{10, ..., 200}
15	ThreadsPerChild	{10, ..., 200}

A.2 Summary of our computational resrouces

All of our data collection experiments were run on a cluster with 20 nodes, each of which is equipped with Intel® Core™ i7-8700 CPU@3.10GHz and 16GB memory. Evaluating all 86M configurations from the 3 systems with 5 repetitions took about 6 months to complete, which results in a total of more than 86,400 CPU hours. For the landscape construction and analyses, all the experiments were carried out using a single node with Intel® Xeon® Platinum 8260 CPU@2.40GHz and 256GB memory.

B Fitness Landscape Analysis

By representing the software configuration landscape as a directed graph⁸, many classic fitness landscape analysis (FLA) methods can be implemented in straightforward graph traversal manners. Here,

⁸Implemented using **NetworkX** package: <https://networkx.org/>.

Table 4. Selected configuration options for SQLite

Index	Parameter	Value
1	SQLITE_SECURE_DELETE	{on, off}
2	SQLITE_TEMP_STORE	{0, 1, 2, 3}
3	SQLITE_ENABLE_AUTO_WRITE	{on, off}
4	SQLITE_ENABLE_STAT3	{on, off}
5	SQLITE_DISABLE_LFS	{on, off}
6	SQLITE_OMIT_AUTO_INDEX	{on, off}
7	SQLITE_OMIT_BETWEEN_OPT	{on, off}
8	SQLITE_OMIT_BTREECOUNT	{on, off}
9	SQLITE_OMIT_LIKE_OPT	{on, off}
10	SQLITE_OMIT_LOOKASIDE	{on, off}
11	SQLITE_OMIT_OR_OPT	{on, off}
12	SQLITE_OMIT_QUICKBALANCE	{on, off}
13	SQLITE_OMIT_SHARED_CACHE	{on, off}
14	CacheSize	{1, ..., 10, 240}
15	AutoVacuumON	{0, 1, 2}
16	ExclusiveLock	{on, off}
17	PageSize	{1, ..., 10, 240}
18	Wal	{on, off}

Table 5. Lookup table of settings of different running environments for three configurable software systems.

Sys.	LLVM	SQLite		Apache	
Index	program_name	num	value_size	requests	concurrency
1	2mm	10	100	50	50
2	3mm	10	1,000	100	100
3	atax	10	10,000	100	100
4	correlation	10	30,000	200	200
5	covariance	100	100	250	250
6	deriche	100	100	300	300
7	doitgen	100	1,000	400	400
8	fdtd2d	100	10,000	500	500
9	gemm	100	30,000	1,000	100
10	symm	1,000	10		
11	syr2k				
12	syrk				
13	trmm				

we delineate the essential ideas and implementations of several FLA methods used in this paper.

Local optima. A local optimum is a configuration that has no superior neighbor. Once the landscape is represented as a graph, the local optima can be easily identified by finding the nodes with no outgoing edges, i.e., the *sink* nodes.

Basin of attraction. While a rugged landscape can be difficult to optimize due to the presence of various local optima, not all are equal in terms of the capability of trapping a solver. For a 2D minimization case, this can be envisioned by the fact that each local optimum is located at the bottom of a ‘basin’ in the landscape surface. Configurations in each basin would eventually fall into the corresponding basin bottom, i.e., the local optimum, when following a simplest hill-climbing local search. To determine the basin of attraction of each local optimum in the landscape, we consider two most popular local search paradigms:

- **Best-improvement local search** (Algorithm 2): In each iteration, the search moves to the neighbor with the highest fitness value. It terminates when no neighbor has a higher fitness value than the current configuration (i.e., a local optimum). For a graph-based landscape, this can be achieved by iteratively selecting the best *successor* of each node until a local optimum

Algorithm 2: Best-Improvement Local Search

Input: A starting configuration \mathbf{c} ; A neighborhood function \mathcal{N} ; A fitness function f

Output: A local optima configuration \mathbf{c}^ℓ

```
1 while  $\mathbf{c}$  is not a local optimum do
2    $\mathbf{c}^* = \operatorname{argmax}_{\mathbf{c}' \in \mathcal{N}(\mathbf{c})} (f(\mathbf{c}'))$  ;
3   if  $f(\mathbf{c}^*) > f(\mathbf{c})$  then
4      $\mathbf{c} \leftarrow \mathbf{c}^*$ ;
5   else
6      $\mathbf{c}$  is a local optimum;
7     break;
```

Algorithm 3: First-Improvement Local Search

Input: A starting configuration \mathbf{c} ; A neighborhood function \mathcal{N} ; A fitness function f

Output: A local optima configuration \mathbf{c}^ℓ

```
1 while  $\mathbf{c}$  is not a local optimum do
2   Improve = False;
3   for  $\mathbf{c}' \in \mathcal{N}(\mathbf{c})$  do
4     if  $f(\mathbf{c}') > f(\mathbf{c})$  then
5        $\mathbf{c} \leftarrow \mathbf{c}'$ ;
6       Improve = True;
7       break;
8   if not Improve then
9      $\mathbf{c}$  is a local optimum;
10    break;
```

is encountered. The *best improvement basin* of a local optimum can be then determined by exhaustively perform such search from each configuration in the landscape, and collect all the configurations that fall into the same local optimum. Note that while this sounds like a computationally expensive task, in practice for landscapes with even millions of configurations, it would take only a few seconds to determine the basin of attraction of each local optimum.

- **First-improvement local search** (Algorithm 3): Here in each iteration, instead of selecting the best neighbor, the search moves to the first neighbor that it encounters with a higher fitness value. This is implemented by iteratively random selecting a successor of each node until a local optimum is reached. Under this paradigm, identifying the basin of attraction of each local optimum is equal to finding the *ancestors* of a node.

Autocorrelation. This is a widely used metric for characterizing the ruggedness of a landscape. As briefly introduced in the main paper, it is the autocorrelation ρ_a of a consecutive series of fitness values $\{f_1, \dots, f_n\}$ obtained from a random walk on the landscape. Due to the graph representation of the landscape, performing a random walk on the landscape is equivalent to that on a graph, which can be executed in a lightning fast manner in **NetworkX**.

Graph embedding. In this paper, we adopted **GL2Vec**, an improved version of **Graph2Vec** to extract low-dimensional features from the LON of each landscape. The generated features are able to capture the topological structure of the LON, and thereby the distribution and connectivity pattern of local optima in the landscape. We employed the implementation of **GL2Vec** from the **Karateclub** package.

Table 6. Full results related to Figure 1 in tabular form, including the number of local optima in each landscape, and the statistics and p -value for comparing the distribution of local optima with random configurations sampled from each landscape.

System	Workload	n peaks	Stat.	p -value
LLVM	gemm	25,974	0.100	$7.6e^{-6}$
	3mm	17,918	0.145	$4.6e^{-2}$
	syrk	43,522	0.095	$2.9e^{-4}$
	trmm	24,658	0.092	$1.7e^{-4}$
	fdtd2d	39,012	0.092	$1.7e^{-4}$
	correlation	12,782	0.145	$4.6e^{-2}$
	2mm	17,494	0.145	$4.6e^{-2}$
	covariance	20,495	0.145	$4.6e^{-2}$
	syr2k	24,619	0.139	$2.5e^{-2}$
	deriche	36,631	0.092	$1.7e^{-4}$
	doitgen	32,148	0.092	$1.7e^{-4}$
	symm	42,962	0.095	$2.9e^{-4}$
	atax	15,388	0.092	$1.7e^{-4}$
APACHE	100_100	204,372	0.067	$9.4e^{-6}$
	200_200	184,939	0.067	$9.4e^{-6}$
	50_50	186,086	0.067	$9.4e^{-6}$
	1000_100	121,153	0.067	$9.4e^{-6}$
	300_300	165,537	0.067	$9.4e^{-6}$
	150_150	179,419	0.067	$9.4e^{-6}$
	250_250	177,716	0.067	$9.4e^{-6}$
	400_400	191,534	0.067	$9.4e^{-6}$
	500_500	128,825	0.067	$9.4e^{-6}$
SQLITE	1000_10	74,261	0.056	$3.2e^{-18}$
	100_100	75,324	0.056	$3.2e^{-18}$
	100_1000	72,908	0.056	$3.2e^{-18}$
	100_10	77,017	0.056	$3.2e^{-18}$
	100_10000	71,053	0.111	$2.9e^{-5}$
	10_1000	79,035	0.111	$2.9e^{-5}$
	10_30000	75,645	0.056	$3.2e^{-18}$
	10_10000	79,438	0.056	$3.2e^{-18}$
	10_100	78,869	0.056	$3.2e^{-18}$
	100_30000	66,680	0.111	$2.9e^{-5}$

C Full Results for Sections 4.1 and 4.2

The full results related to Sections 4.1 to 4.2 can be found in Table 6, Fig. 8, and Fig. 9.

D Full Results for Sections 4.3 and 4.4

The full results related to Sections 4.3 to 4.4 can be found in Figures 10 to 12.

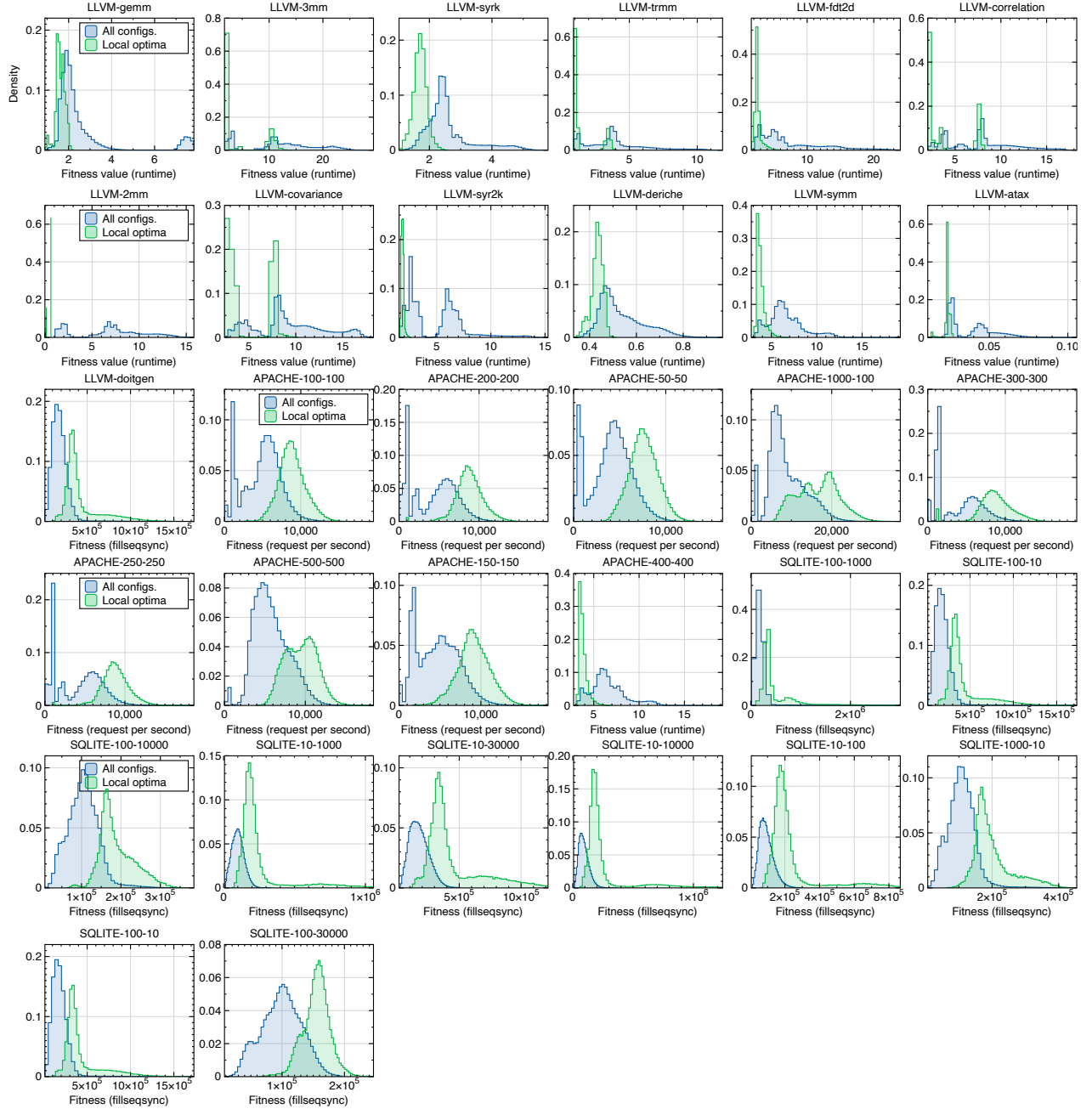


Figure 8. Comparison of fitness distribution of all configurations (blue) versus local optima (green) for all studied landscapes. Note that while the objective function for LLVM is minimized, APACHE and SQLITE have maximized objectives.

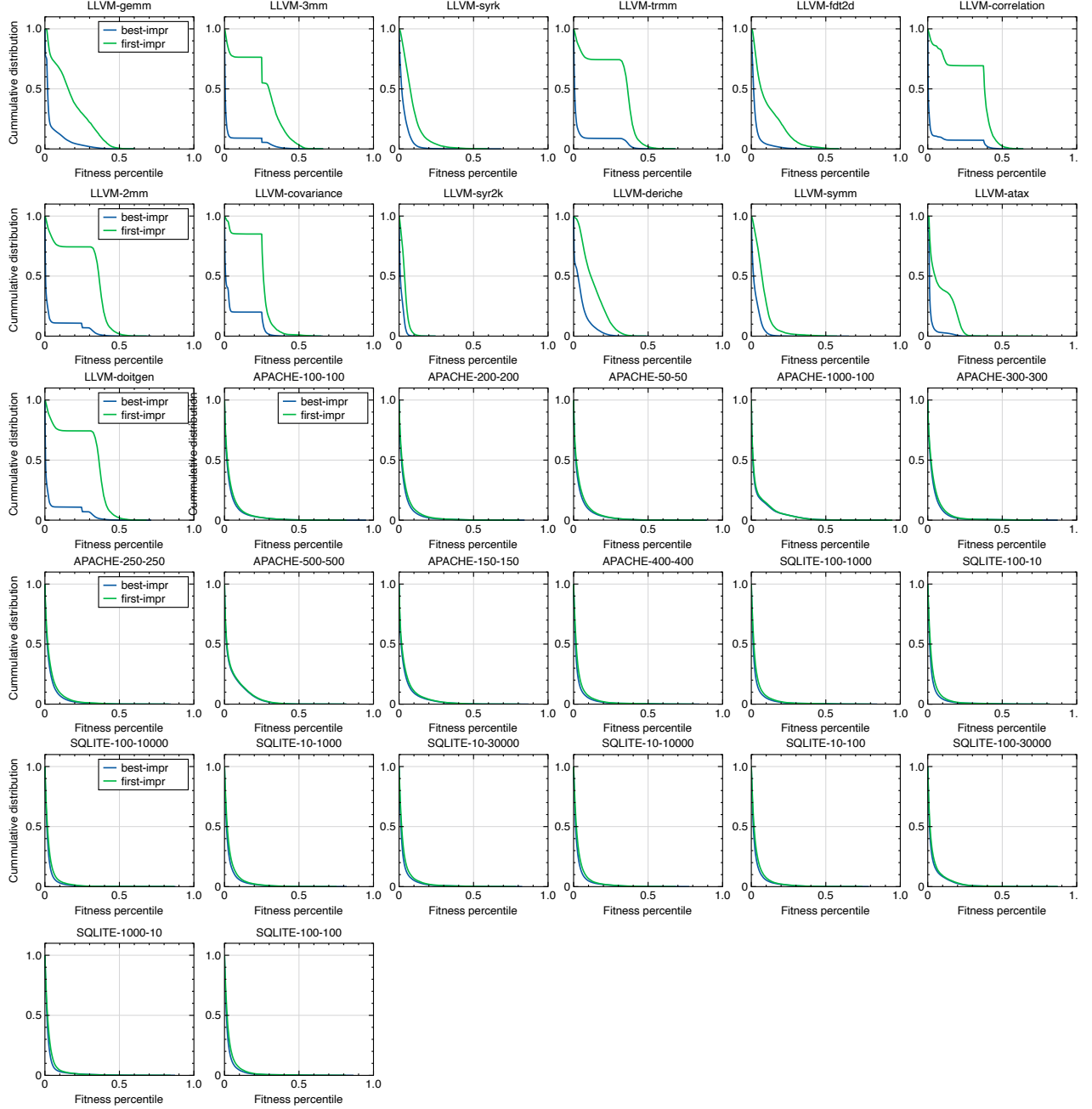


Figure 9. This plot provides the result for the Figure 3(A) and (B) (in the main text) across all scenarios, showing the cumulative distribution of basin size versus the fitness percentile of local optima under both best- and first-improvement local search. Note that for best-improvement, the basin size can be deterministically calculated, and $y = 1.0$ indicates the total sum of all basin sizes, which equals to the total number of configurations in the landscape. For first-improvement, the curves are approximated by conducting 10^9 runs of randomized local search on the landscape, and hence $y = 1.0$ represents the total frequency of visits (i.e., 10^9).



Figure 10. Evolutionary trajectories of warm-start GA (blue) against its vanilla version (purple) in 32 workloads. In particular, both algorithms are started with an initialized population of 50 and the total number of function evaluations is set to 5,000. From these trajectories, we can see that the warm-start GA outperforms its vanilla counterpart, in terms of approximated optimal solution and the convergence rate, in over 78% cases.

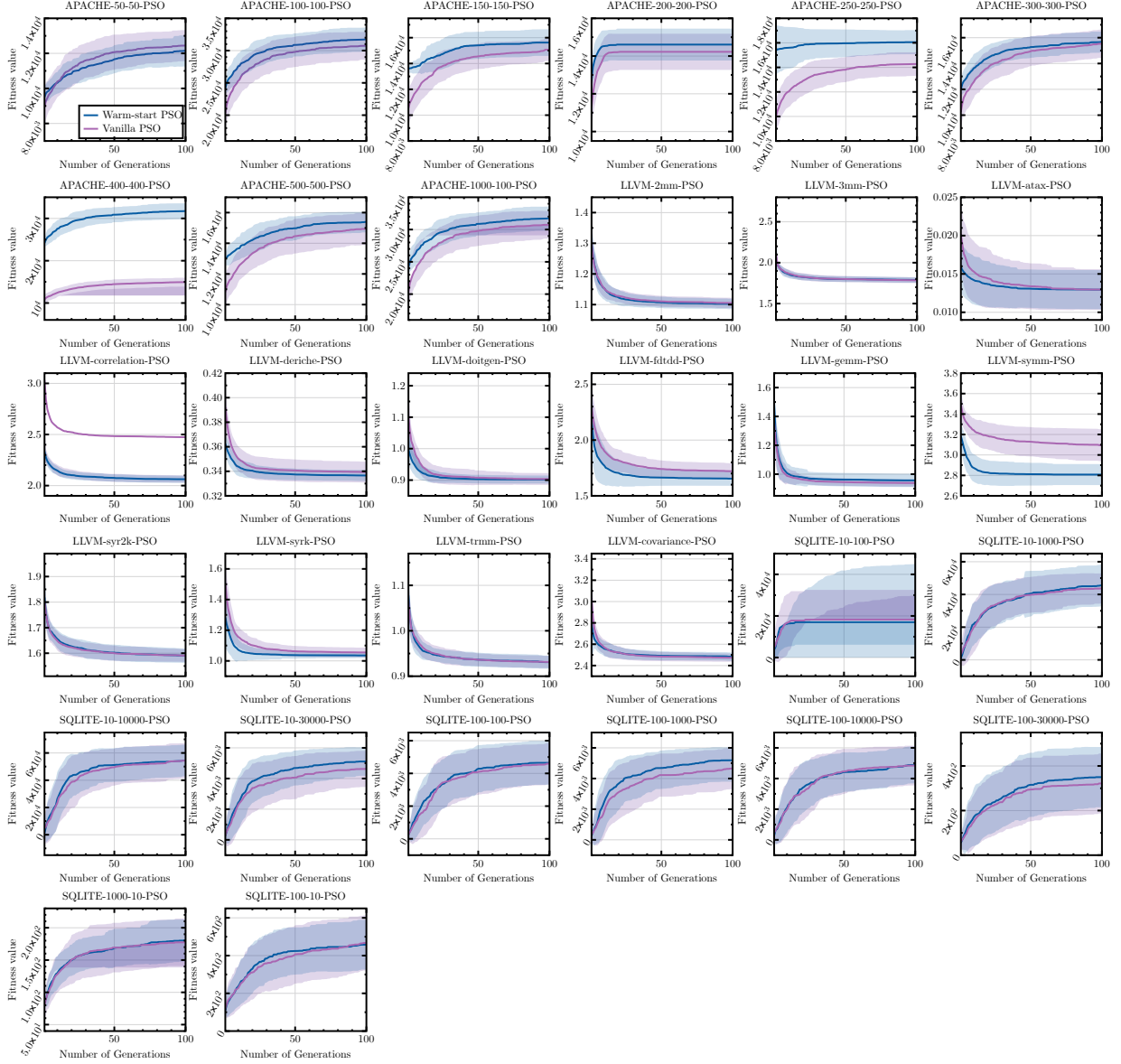


Figure 11. Evolutionary trajectories of warm-start PSO (blue) against its vanilla version (purple) in 32 workloads. In particular, both algorithms are started with an initialized population of 50 and the total number of function evaluations is set to 5,000. From these trajectories, we can see that the warm-start PSO outperforms its vanilla counterpart, in terms of approximated optimal solution and the convergence rate, in 75% cases.

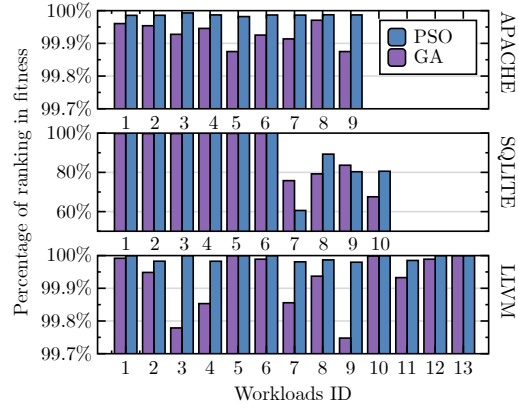


Figure 12. Bar charts of the average percentile of local optima approached by GA and PSO for LLVM, APACHE, and SQLITE on different workloads. From these comparison results, we can see that the solutions obtained by both GA and PSO are located in the top 0.1% local optima whose fitness value exceeds the 99.9% other local optima.