*Original Article*

# Resource Allocation Optimization in Public Cloud Using Genetic Algorithms: A Performance-Cost Trade-off Analysis

**Khaja Kamaluddin**

*Masters in Sciences, Fairleigh Dickinson University, Teaneck, NJ, USA, Aonsoft International Inc, 1600 Golf Rd, Suite 1270, Rolling Meadows, Illinois, 60008 USA.*

**Abstract:** *Resource allocation in public cloud environments requires balancing performance metrics such as execution time and SLA compliance with cost-related factors like resource pricing and utilization. This review examines the effectiveness of Genetic Algorithms (GAs) in addressing this performance–cost trade-off. Gases, with their population-based evolutionary approach, are shown to be well-suited for multi-objective optimization, offering diverse, near-optimal solutions under dynamic cloud workloads. The article discusses key modeling strategies, fitness functions, and the role of GA variants such as NSGA-II in generating Pareto-optimal allocations. A comparative analysis with other metaheuristics, including PSO and ACO, highlights GA's advantages in scalability, adaptability, and trade-off management. The review concludes by identifying research gaps and proposing directions for future work, including hybrid and real-time adaptive GA-based cloud schedulers.*

**Keywords:** *Resource Allocation, Optimization, Cloud Computing.*

## I. INTRODUCTION

Over the past decade, cloud computing has revolutionized the way computing resources are accessed, managed, and delivered. Among the three service models (IaaS, PaaS, and SaaS), Infrastructure-as-a-Service (IaaS) stands out for offering flexible, scalable, and on-demand access to computing infrastructure such as virtual machines, storage, and networks [1]. Public cloud platforms such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP) have transformed how organizations access and scale computing resources. Instead of investing heavily in physical infrastructure, businesses can now deploy applications, store data, and expand their operations on-demand with just a few clicks. This ease of access and scalability, however, introduces new complexities particularly when it comes to efficiently managing and allocating those resources. In practice, cloud resource allocation means deciding how best to assign computing power, storage, and network capacity to different users or applications, in a way that not only ensures strong performance but also keeps operational costs under control.

Despite being a fundamental function, efficient resource allocation is inherently complex due to dynamic workloads, heterogeneous resource requirements, and varying Service Level Agreements (SLAs) as improper allocation can lead to issues such as resource underutilization, SLA violations, increased energy consumption, or inflated operational costs [2]. These challenges become even more complicated in multi-tenant environments, where multiple users share the same cloud infrastructure. In such cases, resources need to be distributed fairly so that no user's performance suffers and at the same time, the overall cost must stay within reasonable limits. Traditional resource allocation approaches in cloud systems include greedy algorithms, heuristics, and rule-based policies. While these methods are effective in simpler scenarios, they struggle with the dynamic and multi-objective nature of real-world cloud environments. Consequently, the research community has turned toward more robust techniques, particularly metaheuristic algorithms, which are better suited for complex, nonlinear, and multi-dimensional optimization problems.

Metaheuristic algorithms like Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO), Simulated Annealing (SA), and Genetic Algorithms (GA) offer powerful, flexible ways to find good solutions to complex problems often in a reasonable amount of time. What makes these techniques especially useful is their ability to adapt to problems of different sizes and types, all while juggling multiple, often conflicting goals like improving performance without driving up costs. Among the metaheuristic approaches, Genetic Algorithms (GAs) have attracted significant interest for their robustness and adaptability in solving complex resource allocation problems in the cloud. Inspired by the principles of natural selection and genetics, GAs simulate evolutionary processes to iteratively improve solutions. Their population-based search mechanism enables exploration of a wide solution space, and their operators (selection, crossover, mutation) help avoid local optima.

One of the key strengths of Genetic Algorithms in cloud resource allocation is their ability to handle multiple goals at once for example, finding the right balance between boosting performance (like improving response time or throughput) and keeping costs down (such as minimizing pricing or energy usage) [3]. Additionally, GAs can adapt to the heterogeneous nature of cloud resources and the dynamic behavior of cloud workloads, making them well-suited for real-time or near real-time scheduling tasks. This article presents a review of the use of Genetic Algorithms for resource allocation in public cloud environments, with a particular focus on the trade-off between performance and cost. While several studies have explored various aspects of GA-based cloud optimization, there is a noticeable gap in literature that offers a consolidated view specifically addressing how GAs manage the performance-cost dilemma in practical scenarios.

The key objectives of this review are:
- To understand how Genetic Algorithms are applied for resource allocation in public cloud computing.
- To classify and compare different approaches, variants, and implementations of GAs used in literature.
- To analyze how existing solutions evaluate and optimize the trade-off between system performance and cost efficiency.
- To identify the limitations of current methods and propose future research directions that can further enhance GA-based solutions.

The scope of the review is limited to academic and industrial research, with emphasis on public cloud platforms. Both single-objective and multi-objective GA models are considered, and their impact on key performance indicators such as execution time, resource utilization, SLA adherence, and cost savings are evaluated. In the following sections, we delve deeper into the working of Genetic Algorithms, their adaptations for cloud environments, challenges of resource allocation, and a comprehensive literature survey to support the performance-cost trade-off analysis. This review aims to serve as a reference point for researchers and practitioners interested in applying evolutionary computing techniques for optimizing public cloud operations.

## II. OVERVIEW OF GENETIC ALGORITHMS

### A. Principles of Genetic Algorithms

Genetic Algorithms (GAs) are a type of evolutionary algorithm inspired by how natural selection and biological evolution work in the real world. Originally introduced by John Holland in the 1970s, GAs start with a group of possible solutions and improve them gradually over multiple generations. With each cycle, they combine and tweak these solutions just like genes mixing and mutating to eventually arrive at an answer that's either optimal or close to it for the problem at hand. The working of a GA is based on the Darwinian principle of "survival of the fittest." Each individual in the population represents a candidate solution, typically encoded as a string (chromosome). The evolution process involves the application of three main operators:
- Selection: This process chooses individuals based on their fitness scores to participate in the generation of new offspring. Techniques like roulette wheel selection, tournament selection, and rank selection are widely used.
- Crossover (Recombination): In this step, two parent solutions combine to produce one or more offspring. Common crossover techniques include one-point, two-point, and uniform crossover.
- Mutation: This introduces random changes in the offspring's genetic code to maintain diversity within the population and prevent premature convergence to local optima.
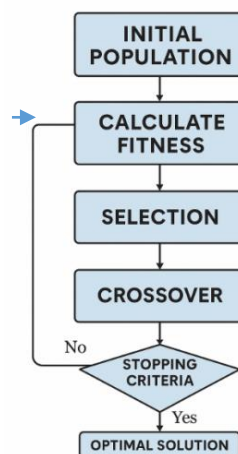


Figure 1: Principles of Genetic Algorithms

The population is evolved iteratively, with each generation expected to improve the fitness of its individuals. The process terminates after a certain number of generations or when a satisfactory solution is found [4].

## B. Advantages of GAs in Cloud Resource Allocation

Genetic Algorithms offer several advantages that make them highly suitable for resource allocation problems in public cloud environments, which are inherently complex, dynamic, and multi-objective:

- Multi-Objective Optimization Capability: GAs can handle multiple conflicting objectives simultaneously, such as minimizing cost while maximizing performance or ensuring load balancing while reducing energy consumption.
- Adaptability to Dynamic Environments: The evolutionary nature of GAs allows them to adapt to changes in the workload, user demands, and resource availability characteristics typical in cloud computing environments.
- Scalability: GAs can be scaled to handle large and complex resource scheduling problems due to their population-based approach.
- Robustness: GAs are less likely to get trapped in local minima compared to traditional heuristics, making them effective for non-linear and high-dimensional problems.
- Simplicity and Flexibility: GAs are relatively easy to implement and can be tailored to a wide variety of optimization problems by designing appropriate fitness functions and encoding schemes.

Due to these characteristics, GAs has become a prominent choice in academic research and practical implementations for cloud resource management.

## C. Common Variants of GAs Used in Literature

Over the years, researchers have proposed various modifications and enhancements to the standard GA framework to better suit specific cloud computing needs. These include:

*a) Multi-objective Genetic Algorithms (MOGAs):*

When an optimization problem has more than one goal like improving performance while also keeping costs low multi-objective Genetic Algorithms (MOGAs) come into play. One of the most popular of these is NSGA-II (Non-dominated Sorting Genetic Algorithm II)**,** which ranks potential solutions based on how well they balance competing objectives and ensures a good variety of options using a technique called *crowding distance*. MOGAs are especially helpful in cloud environments, where it's important to find the right balance between performance and cost, rather than focusing on just one at the expense of the other [5].

*b) Hybrid Genetic Algorithms:*

Hybrid GAs combines the strengths of GAs with other algorithms to improve performance. For example:

- GA + Local Search: Enhances exploitation capability by fine-tuning offspring after crossover.
- GA + Particle Swarm Optimization (PSO): Combines GA's global search with PSO's fast convergence.
- GA + Fuzzy Logic or Neural Networks: Used for intelligent decision-making in resource selection.

Such hybrid techniques have shown promising results in improving scheduling efficiency and reducing response time and energy costs.

*c) Parallel and Distributed GAs:*

To improve scalability and execution speed, parallel GAs distribute the population across multiple processing units. This is especially relevant in cloud environments, where GAs can be parallelized over multiple virtual machines or containers to handle large-scale tasks in real-time.

*d) Elitist and Adaptive GAs*

Elitist GAs preserves the best individuals across generations, ensuring that the highest-quality solutions are not lost. Adaptive GAs dynamically adjusts mutation and crossover rates based on performance trends to avoid stagnation and enhance convergence speed [6].

*e) Domain-Specific Encodings and Operators*

Researchers have also tailored the encoding schemes and genetic operators to better represent cloud-specific parameters. For instance, a chromosome might encode a mapping of tasks to VMs, and the fitness function may consider SLA compliance, cost, and execution time.

## D. Challenges of Applying GAs in the Cloud Context

Despite their strengths, applying GAs to cloud resource allocation also presents several challenges:

- Fitness Function Design: Crafting an effective fitness function that accurately represents cloud objectives (cost, performance, energy, SLA, etc.) is non-trivial and highly domain-specific.

- Execution Time Overhead: In real-time or latency-sensitive environments, the computational cost of GAs may be prohibitive unless parallelized or accelerated.
- Parameter Tuning: The performance of GAs heavily depends on parameters like population size, mutation rate, and crossover rate, which often require empirical tuning.

Nevertheless, these limitations can be mitigated through thoughtful algorithm design, hybridization, and deployment in distributed computing environments. In the next section, we will explore the specific requirements and challenges associated with resource allocation in public cloud environments, setting the stage for understanding how Genetic Algorithms can be tailored to address these complexities.

### III. CLOUD RESOURCE ALLOCATION: REQUIREMENTS AND CHALLENGES

Efficient resource allocation is the backbone of performance and cost management in public cloud computing. With the rapid adoption of cloud-based services, cloud providers must manage a massive pool of heterogeneous resources while meeting the diverse and dynamic needs of clients. This section provides an overview of the types of cloud resources typically involved, outlines key requirements for effective allocation, and highlights major challenges that complicate optimization in real-world cloud environments.

#### A. Resource Types in Public Clouds

Public cloud platforms such as AWS, Azure, and Google Cloud offer a range of resources that can be provisioned on demand. These include:

- Compute Resources: Virtual machines (VMs), containers, and serverless functions that perform processing tasks. These resources vary in CPU architecture, memory, and processing power.
- Memory and Storage: RAM for temporary data handling and storage options like block storage, object storage, and databases (e.g., Amazon S3, Azure Blob).
- Network Resources: Bandwidth, IP addresses, and routing services to ensure seamless communication between cloud services.
- Energy Resources (Indirectly): Although not exposed directly to users, energy consumption is a concern for providers managing data center efficiency.

Each of these resources comes with distinct pricing models, usage policies, and performance characteristics. Resource allocation must therefore not only satisfy the application's functional requirements but also optimize for cost-effectiveness and efficiency.

#### B. Key Requirements for Efficient Resource Allocation

Resource allocation in public cloud computing must meet the following critical requirements:

*a) Performance Assurance*

Applications hosted in the cloud often have strict performance requirements. These are measured in terms of:

- Response time
- Throughput
- Latency
- Execution deadlines

Resource allocation algorithms must ensure that enough resources are allocated to meet these performance targets, even during peak loads.

*b) Cost Efficiency*

Cost is one of the most critical factors for both cloud consumers and providers. The goal is to minimize the total cost of resource usage, which involves:

- Selecting the most cost-effective resource type (e.g., spot vs. on-demand instances)
- Reducing idle or underutilized resources

*c) Scalability and Elasticity*

Cloud systems are expected to scale up or down based on demand. Resource allocation must support:

- Horizontal scaling (adding/removing instances)
- Vertical scaling (changing resource capacity of instances)
- Real-time provisioning and de-provisioning of resources
- Avoiding over-provisioning

*d) SLA Compliance*

Service Level Agreements (SLAs) define performance guarantees between cloud providers and customers. Resource allocation must ensure:

- Minimal SLA violations
- Real-time response to performance degradation
- Penalty avoidance for under-performance

*e) Resource Utilization*

High utilization improves efficiency and reduces waste. Allocation strategies must maximize:

- CPU and memory usage
- Network and I/O utilization
- Storage optimization

*f) Energy Efficiency*

Though more relevant to cloud providers, energy consumption is becoming increasingly important from both cost and environmental perspectives. Green cloud computing aims to reduce energy usage through intelligent resource placement and workload consolidation.



**Figure 2: Requirements for Efficient Resource Allocation in Cloud**

## C. Key Challenges in Resource Allocation

Despite technological advancements, several challenges hinder optimal resource allocation in public clouds [7]:

*a) Dynamic and Unpredictable Workloads*

Cloud workloads are often unpredictable due to varying user demands, usage spikes, and application behavior. Static allocation strategies fail in such scenarios. Real-time or predictive allocation is necessary but hard to implement accurately.

*b) Heterogeneity of Resources and Applications*

Public clouds consist of a mix of hardware and virtualized resources, each with different performance characteristics. Additionally, hosted applications range from web servers to AI training workloads, each with distinct resource needs. This heterogeneity complicates the creation of one-size-fits-all allocation strategies.

*c) Multi-Tenancy and Fairness*

In public clouds, multiple users share the same physical infrastructure. Ensuring fair and isolated resource allocation among tenants while maintaining system-wide efficiency is a complex balancing act.

*d) Trade-off between Cost and Performance*

Maximizing performance typically involves allocating more or premium resources, which increases cost. Conversely, minimizing cost may risk SLA violations or performance degradation. Designing allocation strategies that intelligently balance this trade-off is one of the most pressing challenges.

*e) Scalability of Allocation Algorithms*

As cloud environments grow in size and complexity, resource allocation algorithms must remain scalable. Traditional heuristics often lack scalability and may not deliver optimal performance in large-scale systems.

*f) Latency and Real-Time Constraints*

Some cloud applications require real-time or near-real-time responsiveness. Allocation algorithms that are computationally expensive (like GAs) may be unsuitable unless optimized or parallelized.

*g) Monitoring and Feedback Delays*

Effective allocation requires real-time system state information. However, delays in monitoring and feedback can cause suboptimal decisions, leading to underutilization or over-provisioning.

## IV. PERFORMANCE VS. COST TRADE-OFF IN PUBLIC CLOUD RESOURCE ALLOCATION

In public cloud environments, organizations face a persistent challenge: how to allocate computing resources in a way that ensures high application performance without incurring unnecessary operational costs. This trade-off lies at the heart of cloud computing economics. While cloud providers offer virtually unlimited scalability and on-demand provisioning, these benefits come at a price. Optimizing resource allocation under conflicting objectives such as performance maximization and cost minimization is, therefore, one of the most important and complex tasks in cloud infrastructure management [8]. This section explores the fundamental nature of this trade-off and discusses how resource allocation strategies including those using Genetic Algorithms (GAs) attempt to balance it.

### A. Understanding the Trade-off

The objectives of performance optimization and cost reduction are frequently at odds. Enhanced performance is typically achieved by provisioning additional or higher-tier resources such as virtual machines (VMs) with greater computational power, GPU acceleration, or high IOPS storage which invariably incur higher costs. On the other hand, minimizing costs often entails using fewer or lower-capacity resources, which can compromise response times, throughput, and user experience.

This inverse relationship creates a multi-objective optimization problem, where no single solution is optimal across all objectives. For instance, deploying an increased number of VMs may improve parallelism and reduce makespan but also elevate the hourly billing amount. Similarly, aggressive auto-scaling during peak loads may prevent SLA breaches but significantly increase short-term expenditure if not dynamically regulated. Therefore, an effective resource allocation strategy must identify configurations that strike an optimal balance between these competing goals.

### B. Performance and Cost Metrics

In the context of cloud computing, performance and cost are quantified using distinct yet interrelated metrics. Performance metrics generally reflect the quality of service experienced by end-users and the efficiency of system resource utilization. Among the most commonly referenced performance indicators is makespan, which represents the total execution time required to complete all scheduled tasks [9]. Closely related is response time, defined as the delay between the submission of a task and the initiation of its processing. Throughput is another critical metric, indicating the number of tasks completed per unit of time, thereby serving as a measure of system productivity. Additionally, resource utilization the extent to which allocated resources such as CPU, memory, or network bandwidth are actively used serves as an indicator of system efficiency.

**Table 2: Performance and Cost Metrics**

| Metric Category | Metric | Description |
|---|---|---|
| Performance | Makespan | Total time required to execute all scheduled tasks. |
| | Response Time | Delay between the submission of a task and the beginning of its execution. |
| | Throughput | Number of tasks completed per unit time, reflecting overall system productivity. |
| | Resource Utilization | Percentage of CPU, memory, or network resources actively used during operation. |
| | SLA Violation Rate | Proportion of requests that fail to meet agreed service-level thresholds or deadlines. |
| Cost | Monetary Cost | Total financial expenditure based on resource usage time and instance types. |
| | Instance Pricing Model | Impact of selecting on-demand, reserved, or spot instances on cost optimization. |
| | Idle Resource Overhead | Cost associated with underutilized or unused resources that are still allocated. |
| | Energy Consumption | Power used by data center infrastructure, relevant to provider costs and green computing goals. |

Lastly, SLA violation rate quantifies the proportion of service requests that fail to meet predefined service level agreement thresholds, which is especially important in commercial cloud settings where penalties may apply for non-

compliance. Cost metrics, on the other hand, evaluate the financial and operational expenditure associated with resource allocation. Monetary cost is the most direct measure, encompassing the cumulative charges incurred based on resource usage duration, instance type, and configuration, as dictated by cloud providers' billing models typically per hour or per second.

The pricing model of instances, such as on-demand, reserved, or spot pricing, plays a substantial role in determining overall cost efficiency. Another important factor is idle resource overhead, referring to the cost of maintaining allocated resources that are underutilized or remain unused for significant periods. While not always visible to cloud tenants, energy consumption represents a backend cost relevant to service providers; its reduction is central to green computing initiatives and indirectly affects long-term pricing and sustainability practices.

In multi-objective optimization problems such as those addressed by Genetic Algorithms these metrics are often either combined into a single composite fitness function using weighted parameters or treated as separate objective functions. The latter approach enables the generation of Pareto-optimal solutions, wherein performance cannot be improved without a corresponding increase in cost, and vice versa.

## C. Multi-Objective Optimization and Trade-off Representation

The performance–cost trade-off can be modeled using single-objective or multi-objective formulations. In single-objective models, a weighted sum approach is typically employed as in (1):

$$Fitness = \alpha \cdot \left(\frac{1}{Makespan}\right) + \beta \cdot \left(\frac{1}{Cost}\right) \qquad (1)$$

Where $\alpha$ and $\beta$ represent the relative importance of performance and cost, respectively. However, such approaches require careful tuning of weights and may bias the search toward one objective.

Multi-objective optimization models, particularly those employing evolutionary techniques such as NSGA-II studied in [10] treat cost and performance as independent objectives. The result is a set of Pareto-optimal solutions, where improvement in one objective necessitates compromise in the other which allows decision-makers to select a suitable trade-off point based on contextual constraints or preferences.

## D. Application of Genetic Algorithms to Trade-off Optimization

The application of Genetic Algorithms (GAs) to cloud resource allocation has emerged as a well-established approach in addressing multi-objective problems, particularly those involving conflicting goals such as performance enhancement and cost reduction. GAs offer a population-based, heuristic-driven search mechanism that excels at identifying near-optimal solutions in large, complex decision spaces. Their evolutionary structure, inspired by the biological processes of selection, crossover, and mutation, allows them to adaptively explore the solution landscape and converge toward configurations that provide an acceptable trade-off between objectives. In resource allocation for public clouds, each candidate solution is typically encoded as a chromosome, with each gene representing a decision variable such as which virtual machine (VM) a task is mapped to, or which type of instance (e.g., spot, on-demand, reserved) is selected for a given workload. This encoding allows for a diverse and flexible representation of the scheduling space, which is particularly important given the heterogeneity of cloud resources and dynamic workload patterns.

The evaluation of solutions in GAs depends on a fitness function that reflects the system's performance and cost objectives. As discussed in previous section, scalarized or normalized fitness formulations are frequently used to evaluate trade-offs. However, in practical implementations, the strength of GAs lies not just in their ability to minimize cost or makespan independently, but in their capability to generate a set of diverse solutions, each reflecting a different point on the performance–cost spectrum. This is especially valuable in cloud environments where contextual needs (e.g., SLA strictness, budget limits, burst load tolerance) vary over time. To efficiently manage such trade-offs, many implementations have adopted multi-objective genetic algorithms (MOGAs) such as NSGA-II and SPEA2. These algorithms maintain a Pareto front of non-dominated solutions, enabling decision-makers or automated schedulers to select configurations dynamically. The work by Singh et al. in [11], for instance, demonstrated that NSGA-II-based resource allocation in simulated cloud platforms (e.g., CloudSim) consistently outperformed baseline heuristics by offering better diversity in trade-off solutions without compromising convergence speed.

A key advantage of GA-based methods is their adaptability to dynamic, real-time environments and as cloud workloads are seldom static they fluctuate based on user demand, application logic, and external triggers such as promotions or system updates. GAs can be adapted to accommodate such dynamics by integrating real-time feedback mechanisms. For instance, adaptive mutation rates, re-evaluation of fitness during execution, and hybrid strategies (e.g., GA + Ant Colony Optimization) have been proposed to ensure timely convergence even in rapidly evolving operational contexts. From an architectural perspective, GA-based resource allocators can be implemented as part of cloud middleware or integrated into orchestrators managing virtual machines, containers, or services. Their computational cost, often cited as a limitation, has been addressed

through parallelization. Studies such as [12], explored GPU-based parallel GAs to reduce computation time for large-scale job scheduling without compromising quality of results. These parallel models are particularly suited to data centers and high-throughput cloud computing environments where scalability is critical.

Moreover, GAs provides flexibility to incorporate user-defined constraints, such as power budgets, geographic availability zones, or compliance rules, which is not always feasible with deterministic or greedy algorithms. This makes them particularly attractive for organizations operating in regulated environments or with sustainability objectives, where trade-offs go beyond just performance and cost.

### E. Implications for Cloud Resource Management

Understanding and modeling the performance–cost trade-off is essential for practical cloud orchestration. Application types differ significantly in their sensitivity to cost and performance. For example interactive applications such as online gaming or real-time analytics prioritize low latency and SLA compliance, justifying higher cost. Secondly, batch processing workloads such as data backups or report generation may tolerate performance degradation in exchange for cost savings. Additionally, resource allocation frameworks that employ GAs are capable of adapting to such distinctions through appropriately designed fitness functions and evolutionary strategies.

### VI. COMPARATIVE ANALYSIS OF GENETIC ALGORITHMS WITH OTHER METAHEURISTICS

The landscape of metaheuristic algorithms for cloud resource allocation is rich and diverse, with Genetic Algorithms (GAs), Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO), and Simulated Annealing (SA) being among the most prominent techniques investigate. While each of these algorithms offers strengths in handling complex optimization problems, their suitability for performance–cost trade-off optimization in public cloud environments varies significantly based on the problem structure, scalability needs, and adaptability requirements. This section presents a comparative analysis of GAs with other widely used metaheuristics in the context of multi-objective scheduling and resource provisioning in Infrastructure-as-a-Service (IaaS) cloud environments.

### A. Comparison with Particle Swarm Optimization (PSO)

PSO, inspired by the social behavior of bird flocking, has been widely applied to scheduling problems due to its simplicity and convergence speed. Each solution (or "particle") adjusts its position in the search space based on its own best performance and that of its neighbors. Studies such as [13] have shown that PSO achieves faster convergence than GAs in homogeneous workload environments. However, PSO tends to suffer from premature convergence and local optima trapping, particularly in highly dynamic or nonlinear scheduling spaces, which are common in public cloud deployments.

In contrast, GAs maintain a more diverse population due to the use of crossover and mutation operators. This diversity preservation allows GAs to explore broader regions of the solution space, often leading to better global optimality in performance–cost trade-offs. For example, experimental evaluations in CloudSim [14] indicated that GA-based schedulers maintained lower SLA violation rates under fluctuating workloads compared to PSO-based alternatives.

### B. Comparison with Ant Colony Optimization (ACO)

ACO simulates the pheromone trail-based pathfinding behavior of ants and has demonstrated strong performance in discrete optimization tasks, particularly task scheduling and network routing. In cloud scheduling problems, ACO has been employed to minimize execution time and maximize resource utilization [15].

**Table 2: Comparative Analysis**

| Feature | Genetic Algorithm (GA) | Particle Swarm Optimization (PSO) | Ant Colony Optimization (ACO) | Simulated Annealing (SA) |
|---|---|---|---|---|
| Population-based Search | Yes | Yes | No | No |
| Exploration vs. Exploitation | Balanced | Exploitation-biased | Exploitation-biased | Local search |
| Multi-objective Optimization | Strong (e.g., NSGA-II) | Moderate (MOPSO) | Limited | Poor |
| Scalability | High | Moderate | Low–Moderate | Low |
| Diversity Preservation | Strong (mutation/crossover) | Weak (leader-focused) | Medium (pheromone balance) | Weak |
| Implementation Complexity | Moderate | Low | High | Low |
| Suitability for Cloud Systems | High | Medium | Medium | Low |

However, ACO generally requires extensive parameter tuning (e.g., pheromone evaporation rate, trail influence), which complicates implementation in real-time cloud systems. GAs, by contrast, are comparatively easier to generalize across problem instances and can naturally accommodate hybrid models. Hybrid GA-ACO frameworks have been proposed to leverage ACO's exploitation strength and GA's exploration ability [16] resulting in superior makespan–cost balances. Nonetheless, standalone GAs remain more stable in performance across varying workload conditions without relying heavily on domain-specific heuristics.

## C. Comparison with Simulated Annealing (SA)

SA is a single-solution based metaheuristic that emulates the cooling process of metals to gradually reduce the acceptance probability of worse solutions. While it can avoid local optima, SA's sequential search nature limits its scalability for large cloud scheduling problems. Its application has largely been confined to small-scale task graphs or static scheduling environments. In contrast, GAs evaluate a population of solutions in parallel, making them inherently more scalable and parallelizable, which is essential for real-world public cloud platforms. Moreover, SA lacks the recombination operator present in GAs, which restricts its capacity to generate novel, diverse solutions a key factor in multi-objective optimization contexts.

## D. Multi-Objective Optimization Capability

A distinguishing advantage of GAs lies in their multi-objective optimization capability, particularly when using Pareto-based extensions such as NSGA-II. While PSO and ACO have been adapted into multi-objective versions (MOPSO, MOACO), these variants often face difficulty maintaining a well-distributed Pareto front due to convergence behavior biased toward local leaders or pheromone trails. GAs, in contrast, can maintain solution diversity through mechanisms such as crowding distance and elitism, ensuring better coverage of the performance–cost trade-off surface. Simulation studies [17] reported that NSGA-II consistently produced more evenly distributed and non-dominated solutions compared to MOPSO in heterogeneous cloud workloads.

## E. Practical Deployability

Deploying metaheuristics in production cloud environments demands not only theoretical efficiency but also robustness to dynamic conditions and ease of integration into orchestration systems. GAs have been shown to perform reliably under uncertain conditions, including changing workload intensity, spot instance interruptions, and SLA fluctuations. Furthermore, their modular architecture allows seamless integration into microservice-based schedulers and container orchestrators (e.g., Kubernetes plugins), especially when implemented in parallel or distributed modes. Although PSO and ACO can also be embedded into scheduling controllers, they often require more careful configuration and do not generalize as well across varying workload types and SLA models. In contrast, GA-based frameworks, when paired with adaptive tuning and feedback-based control loops, can achieve near-real-time resource allocation with balanced cost and performance efficiency.
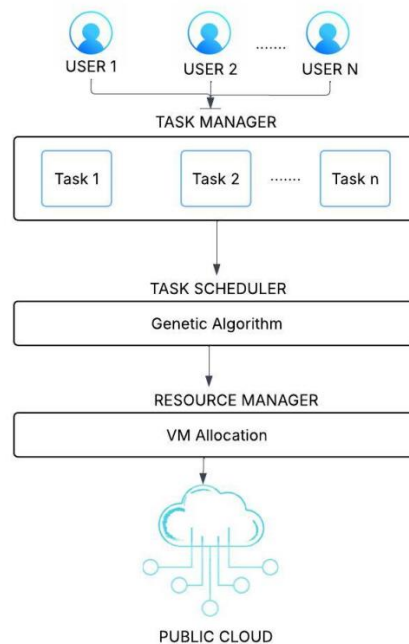


**Figure 3: GA Based Resource Scheduling in Cloud**

Additionally, GA-driven scheduling decisions can be extended to incorporate security-aware resource provisioning, enabling proactive mitigation of cloud-native threats such as co-residency risks and lateral movement during development-stage deployments [18].

## VI. CONCLUSION AND FUTURE DIRECTIONS

Through structured comparisons, we observed that GAs are particularly well-suited for balancing performance metrics such as makespan, response time, and SLA adherence against cost-related concerns including resource billing models, idle overhead, and energy consumption. Their ability to evolve diverse, near-optimal solutions enables the discovery of Pareto fronts, offering decision-makers a set of allocation strategies that can be dynamically selected based on context-specific constraints. Despite their advantages, GA-based approaches are not without limitations. They are computationally intensive and require careful parameter tuning (e.g., population size, mutation rate) to ensure efficient convergence. In real-time environments, the latency introduced by iterative evolution may limit their immediate applicability unless supported by parallel processing or approximation techniques.

Furthermore, while many studies utilize simulated environments such as CloudSim, few have validated GA-based scheduling in production-grade or hybrid/multi-cloud infrastructures. Moving forward, several avenues for future research emerge. First, the integration of GAs with predictive models such as demand forecasting using machine learning could enhance proactive resource planning. Second, real-time adaptive GAs, which dynamically adjust operators based on feedback from cloud performance metrics, offer promising directions for autonomous cloud management. Third, hybrid algorithms that combine GAs with faster local search or swarm-based techniques may help reduce execution time while preserving solution diversity. Lastly, with the rise of serverless computing and edge–cloud continuum models, the application of GAs to decentralized and event-driven environments presents a novel and underexplored frontier.

## REFERENCES

[1] Manvi, S. S., & Shyam, G. K. (2014). Resource management for Infrastructure as a Service (IaaS) in cloud computing: A survey. *Journal of network and computer applications*, *41*, 424-440.

[2] Odun-Ayo, I., Udemezue, B., & Kilanko, A. (2019). Cloud service level agreements and resource management. Adv. Sci. Technol. Eng. Syst, 4(2), 228-236.

[3] Alkayal, E. (2018). Optimizing resource allocation using multi-objective particle swarm optimization in cloud computing systems (Doctoral dissertation, University of Southampton). K. Elissa, "Title of paper if known," unpublished.

[4] Immanuel, S. D., & Chakraborty, U. K. (2019, July). Genetic algorithm: an approach on optimization. In 2019 international conference on communication and electronics systems (ICCES) (pp. 701-708). IEEE.

[5] Selçuklu, S. B. (2023). Multi-objective genetic algorithms. In Handbook of formal optimization (pp. 1-37). Singapore: Springer Nature Singapore.

[6] Ahn, C. W., & Ramakrishna, R. S. (2003). Elitism-based compact genetic algorithms. IEEE Transactions on Evolutionary Computation, 7(4), 367-385.

[7] Abid, A., Manzoor, M. F., Farooq, M. S., Farooq, U., & Hussain, M. (2020). Challenges and issues of resource allocation techniques in cloud computing. KSII Transactions on Internet and Information Systems (TIIS), 14(7), 2815-2839.

[8] Fé, I., Matos, R., Dantas, J., Melo, C., Nguyen, T. A., Min, D., ... & Maciel, P. R. M. (2022). Performance-cost trade-off in auto-scaling mechanisms for cloud computing. Sensors, 22(3), 1221.

[9] Kaur, R., Laxmi, V., & Balkrishan. (2022). Performance evaluation of task scheduling algorithms in virtual cloud environment to minimize makespan. International Journal of Information Technology, 1-15.

[10] Deb, K., Pratap, A., Agarwal, S., & Meyarivan, T. A. M. T. (2002). A fast and elitist multiobjective genetic algorithm: NSGA-II. IEEE transactions on evolutionary computation, 6(2), 182-197.

[11] Singh, M. K., Choudhary, A., Gulia, S., & Verma, A. (2023). Multi-objective NSGA-II optimization framework for UAV path planning in an UAV-assisted WSN. The Journal of Supercomputing, 79(1), 832-866.

[12] Cheng, J. R., & Gen, M. (2020). Parallel genetic algorithms with GPU computing. In Industry 4.0-Impact on Intelligent Logistics and Manufacturing. IntechOpen.

[13] Al Reshan, M. S., Syed, D., Islam, N., Shaikh, A., Hamdi, M., Elmagzoub, M. A., ... & Talpur, K. H. (2023). A fast converging and globally optimized approach for load balancing in cloud computing. IEEE Access, 11, 11390-11404.

[14] Masdari, M., Salehi, F., Jalali, M., & Bidaki, M. (2017). A survey of PSO-based scheduling algorithms in cloud computing. Journal of Network and Systems Management, 25(1), 122-158.

[15] Sharma, N., & Garg, P. (2022). Ant colony based optimization model for QoS-Based task scheduling in cloud computing environment. Measurement: Sensors, 24, 100531.

[16] Zukhri, Z., & Paputungan, I. V. (2013). A hybrid optimization algorithm based on genetic algorithm and ant colony optimization. International Journal of Artificial Intelligence & Applications, 4(5), 63-75.[13

[17] Ding, S., Chen, C., Xin, B., & Pardalos, P. M. (2018). A bi-objective load balancing model in a distributed simulation system using NSGA-II and MOPSO approaches. Applied soft computing, 63, 249-267.

[18] Yashu, M. S. F. (2021). Thread mitigation in cloud native application Develop-Ment.