

Citation: Swathi M, Krishna Sameera V, Satya Prakash K V S, et al. Dynamic load balancing in multi-cloud environments using genetic and ant colony optimization approach: A hybrid approach. *Journal of Harbin Institute of Technology (New Series)*. DOI: 10.11916/j.issn.1005-9113.2025059

Dynamic Load Balancing in Multi-Cloud Environments Using Genetic and Ant Colony Optimization Approach: A Hybrid Approach

Swathi M¹, Krishna Sameera V¹, Satya Prakash K V S², Kiran Kumar B^{3*}, Prathyusha Dogga⁴ and Mamatha Vayalapelli⁵

1. Department of Computer Science and Engineering, Anil Neerukonda Institute of Technology & Sciences, Vishakapatnam 530048, India;
2. Department of Information Technology, Gayatri Vidya Parishad College of Engineering, Vishakapatnam 530048, India;
3. School of Computing, SRM Institute of Science and Technology, Tiruchirappalli 621105, India;
4. Department of Mechanical Engineering, Vignan's Institute of Information Technology, Visakhapatnam 530049, India;
5. Department of Computer Science and Engineering, ARKA Jain University, Jamshedpur 832108, India.)

Abstract: The multiplicity of the clouds improves the availability and elasticity, however, the load balancing is more complex when the resources are not balanced, the loads are not constant, and cross-providers failures occur. The present work suggests the initial dynamical load balancing in real-time and multi-cloud with the aim to optimize the resources of the system, reduce the latency (tail included), and enhance the fault tolerance. This paper presents an account of a hybrid algorithm which is a combination of Genetic Algorithm (GA) and Ant Colony Optimization (ACO). Although GA does not only consider the global search through candidate task-resource mappings, ACO concentrates on local search, using the pheromone search. This publication defines explicit parameter tuning and workflow to be used in evaluation. The hybrid GA-ACO was introduced and was compared to the older ones (Round – Robin, Least-Connections), adaptive methods and ML-based methods in the context of various multi-cloud cases (steady, spike, node-fail, regional-outage). The assessment was done using the real deployment-related KPIs: load distribution, mean/P99 latency, resource utilization, fault tolerance. The hybrid GA-ACO improved its load imbalance by 36%, mean response time by 48%, resource utilization by 37% and fault tolerance by 36%. It is also possible to maintain 150 ms requested binding per control window and less than 5 ms overhead of the request-path. Even though our hybrid solution will reduce imbalance and latency and increase utilization and fault tolerance, the usability and fault tolerance will be a priority in future efforts through extensive validation with real workloads. The proposed hybrid solution is a combination of GA and ACO and represents global and local search using parameter optimization and the outcome of a thorough consideration, offering a smooth, adaptive, and efficient way of load balancing between various clouds.

Keywords: load balancing; genetic algorithm; ACO; multi-cloud; hybrid approach

CLC number: TP301.6

Document code: A

Article ID: 1005-9113(2026)00-0000-16

0 Introduction

The rapid development of the cloud computing has transformed the way in which organizations deploy, operate and scale their IT resources. multi-cloud environments, which are defined as the process of dedication of the services to different cloud providers, have become popular with the aim of enhanced reliability, flexibility, and cost savings^[1]. Though there are numerous benefits of multi-cloud environments, they have serious problems in

managing resources particularly in load balancing^[2]. Cloud computing must involve load balancing, which enables the workload to be distributed between available resources in the most efficient way to achieve better performance, fewer latencies, and higher fault tolerance. Though effective in single-cloud environment, traditional load balancing systems often have to advance in the complex and dynamic environment of the multi-cloud environment^[3]. The irregularity of workload trends, dissimilarities in resource abilities and the need to adjust to circumstances in real-time requires more sophisticated

solutions. This study researches on the emerging algorithm, which incorporates Genetic Algorithms (GA) and Ant Colony Optimization (ACO) to assist in the issues of load balancing across a number of clouds^[4]. GA are also known to be able to do global search which enables a wide range of possible solutions to be explored. GA mimic the natural selection to find the best solutions to complex optimization problems in a very efficient manner^[5-6].

On the other hand, ACO proves to be better in local search optimization. ACO is a search algorithm based on the foraging behaviour of ants, which uses pheromone trails for solution direction, and optimizes the solution based on cooperative mechanisms^[7-8]. The suggested hybrid GA-ACO algorithm aims at taking the benefits of both optimization methods. The GA component will explore different load distribution plans on the international scale, however, the ACO component will optimize the load distribution plans to achieve optimum performance at the local level. This two-fold strategy is designed to maximize resource usage, reduce the response time, and increase the reliability of systems by dynamically adapting to the changes in demand that are seen in the multi-cloud environment. Extensive simulation and experimental work will be done to prove the efficiency of the hybrid GA-ACO algorithm. Key metrics will be used to assess the performance of the algorithm against the existing load balancing techniques and they include the efficiency of distributing loads, response time, resource consumption, and fault tolerance.

This study is expected to produce useful findings and make an important contribution to the development of load-balancing techniques in multi-cloud computing. The goal of this research is to develop a strong and dynamic load-balancing platform that will specifically be created to support multi-cloud configuration. The proposed hybrid algorithm will be based on the outstanding characteristics of genetic algorithm as a global optimization algorithm and the ant colony algorithm as a local optimization algorithm. It addresses the issues of multi-cloud load balancing, and it accelerates and becomes more dependable cloud computing infrastructures.

Problem statement: Consider a multi-cloud environment with regions/resources R and compute instances/VMs $M = \cup_{r \in R} M_r$. A set of incoming tasks/jobs T must be assigned dynamically to resources.

R denotes the set of regions or cloud resource domains in the multi-cloud environment. M represents the global set of compute instances or Virtual Machines (VMs), defined as the union of all regional resource sets, i.e., $M = \cup_{r \in R} M_r$, where M_r is the set of compute instances available in region $r \in R_r$. T denotes the set of incoming tasks or jobs that must be dynamically scheduled. The objective is to assign each task in T to an appropriate compute instance in M across regions. Here, each M_r represents the set of VMs available in a specific region r , and the union operation aggregates all regional VM sets into a single global resource pool M .

Parameters:

τ_i : SLA deadline/latency target for task i .

B_r : optional bandwidth limit for region r ;

b_{ij} : bandwidth consumed if $i \rightarrow j$.

Decision variables $x_{ij} \in \{0,1\}$: 1 if task i is assigned to VM_j ; 0 otherwise. **Per-task latency model:** Let service time contribution on j be $S_{ij} = d_i/c_j$; **Predicted end-to-end latency:** $T_{ij} = L_{ij} + S_{ij}$. Here, c_j represents effective compute capacity of VM_j (tasks/s or normalized units). L_{ij} : network/queuing latency for assigning task i to VM_j (ms).

Objective function:

$$\min_{x,y} w_{lat} \sum_{i \in T} \sum_{j \in M} x_{ij} T_{ij} + w_{sla} \sum_{i \in T} \sum_{j \in M} x_{ij} [T_{ij} - \tau_i] + w_{cost} \sum_{i \in T} \sum_{j \in M} x_{ij} (k_j d_i + \varepsilon_{ij}),$$

where $[\mu]_+ = \max\{\mu, 0\}$ penalize SLA violations. κ_j : per-unit compute cost at VM_j ; d_i : compute demand of task i (normalized service units). ε_{ij} : egress/data-transfer cost if $i \rightarrow j$. $w_{lat}, w_{sla}, w_{cost} \geq 0$: objective weights (set by operator).

This paper makes the following contributions.

1) **Hybrid dual-timescale optimizer:** A GA-ACO design where GA performs inter-cloud placement/share allocation and ACO performs intra-cloud task-to-VM routing, enabling coordinated decisions across clouds and within regions.

2) **Joint latency-cost-SLA objective with tail control:** A multi-objective fitness incorporating response time, SLA violations, and egress/compute cost, augmented with a risk-sensitive (CVaR (Conditional Value at Risk)/quantile) term to limit high-percentile latency.

3) **Robustness to shift:** A generalization protocol using domain randomization and scenario-level cross-validation (no leakage), plus regularizes (fitness noise, early stopping, diversity maintenance) to

prevent overfitting to simulation scenarios.

4) Automated + online parameter control: A two-stage scheme, offline BOHB (Bayesian Optimization with HyperBand)/Bayesian auto-tuning to find robust settings, and online adaptive rules (driven by diversity, SLA pressure, load variance) to remain effective under non-stationary workloads.

5) Scalability mechanisms: Island-model GA and multi-colony ACO with asynchronous merges, hierarchical regional decomposition, and incremental (δ) evaluation, reducing wall-clock time and memory pressure on large instances.

1 Related Work

Cloud computing has advanced quite well in load balancing with many approaches developed to address the problem of dynamic and distributed systems. This section looks into the past studies and outlines the contribution and weakness of present approaches in multi-cloud load balancing.

1.1 Conventional Load Balancing Methods

Preliminary studies in load balancing mainly focused on the static algorithms, such as Round Robin, Least Connections, and Weighted Balancing of Loads. These plans apportion tasks as per the laid down standards or measures. Their performance is good in intra-cloud environments; however, they often do not scale to the dynamism of multi-cloud networked environments^[9-10]. These solutions should be more flexible in order to handle varying workloads and availability of resources in different cloud platforms.

1.2 Dynamic and Adaptive Load Distribution

The load-balancing strategies are provided to be dynamic and adaptive in an effort to balance the limitations of the systems which are fixed. A majority of the smaller response time and adaptive load balancing adapt the resource allocation according to the real time measurements and enhance the performance and responsiveness^[11]. Nevertheless, these solutions tend to be based on the centralized control, which may be considered as a drawback in the context of the multi-clouds use where the resources are distributed.

1.3 Approaches Based on Machine Learning

Machine learning algorithm has been used in load balancing to improve the decision making. The use of the reinforcement learning, neural networks, and

other machine learning frameworks in predicting the workloads and optimizing resource allocation has been proposed in the recent studies^[12]. Despite promising outcomes of these methodologies, most of the times they have large training data requirements and are resource-consuming, thus they cannot scale and cannot be deployed in real-time.

1.4 Metaheuristic Optimization Methods

GA and ACO are also metahyperoptical algorithms that have been applied to complex load-balancing problems. According to natural selection, genetic algorithms have been adopted to discover the most optimal or near optimal ways of distributing load by allowing a population of candidate solutions to evolve^[13]. It has also proven that the application of ant colony optimization, which is inspired by the foraging behaviour of ants, to solve routing and scheduling problems is effective in optimization of local searches^[14]. However, there has been a limited use of such methods in multi-cloud load balancing particularly in the integration of their benefits.

1.5 Hybrid Methodologies

More recent research has examined the case of hybrid methodologies which combine multiple optimization methods to take advantage of the benefits. Integration of GA and Particle Swarm Optimization (PSO) or Simulated Annealing (SA) has been demonstrated to enhance performance in a number of optimization problems^[15-16]. A combination of GA with the ACO could help in load balancing, as the global search skills are combined with the local optimization capabilities. The use of hybrid methods in multi-cloud load balancing remains an active research topic.

1.6 Load Balancing Utilizing Blockchain Technology

The blockchain technology is proposed to enhance the process of load balancing through decentralized and transparent decision-making processes^[17]. The application of blockchain based solutions has advantages in terms of security and trust but is not developed in the field of multi-cloud systems dynamic and real-time load balancing. Combining the blockchain technology with optimization methods such as genetic techniques (GA) and ant colony optimization (ACO) reflects a novel area of research^[18]. Existing studies have developed load-balancing policies in cloud systems yet there are still challenges in addressing the dynamic and

distributed nature of multi-cloud systems^[19–20]. Traditional methods, despite being essential, often are more flexible to the modern multi-cloud. The higher order methods, including machine learning and metaheuristic optimization, offer fascinating solutions but they need to be explored to work out all hybrid methodologies. The purpose of this paper is to make a progressive step forward in order to create a hybrid approach to combine genetic algorithm and ant colony optimization in multi-cloud load balancing through incorporating the advantages of both algorithms in order to increase efficiency and flexibility in complicated cloud systems.

1.7 Research Gaps and Positioning

Though many studies have been conducted on the load balancing in both cloud and multi-cloud environment, there are still many research gaps that have not been filled, as shown in Table 1.

To begin with, the majority of current approaches are single-objective, that is, they do not consider multi-objective trade-offs, including SLA violations and fairness. Second, most of them use one metaheuristic, e.g. GA or ACO only, and thus they tend to converge slowly or come to a premature halt. Third, these methods can be manually tuned to

parameters, which limits flexibility with dynamic workloads. Fourth, aspects of scalability and real time deployments are not commonly tested when there is a large scale multi-cloud. Lastly, strength and ability to work on unseen workload patterns is still not studied sufficiently. In order to overcome these loopholes, the current research paper suggests a hybrid Genetic Algorithms-Ant Colony Optimization (GA-ACO) model that combines global exploration and local exploitation in a multi-objective dynamically adaptive model. The suggested model incorporates online adaptive-parameter control, hierarchical decomposition to ensure scalability, and cross-scenario testing to ensure robustness, therefore, it overcomes the critical limitations of the current studies and offers a more adaptive, scalable, and SLA-aware load balancing service.

2 The Proposed Methodology

This section will discuss the mechanism of the hybridization and experimentation of hybrid GA and ACO in the multi-cloud systems in terms of dynamic load balancing. It has three major phases, i. e., creation of the algorithm, implementation and testing, and evaluation of its performance.

Table 1 Differences between existing work and proposed methodology

Aspect	This paper	Typical existing approaches
Decision structure	Dual-timescale GA-ACO; GA for inter-cloud shares; ACO for intra-cloud routing	Single-layer heuristic or single metaheuristic
Objective	Latency + SLA violations + cost with tail-aware (CVaR/quantile) control	Average latency or cost; limited tail control
Adaptively	Online parameter control (diversity/SLA/load signals); portfolio seeding	Mostly fixed parameters; manual re-tuning
Robustness	Domain randomization + scenario-level CV; fitness noise; early stopping	Random-seed CV; limited shift testing
Scalability	Island GA, multi-colony ACO, hierarchical decomposition, delta evaluation	Centralized/monolithic solvers
Deployment path	Kubernetes/Terraform testbed plan; CloudLab/OpenStack prototype	Primarily simulation

2.1 Algorithm Design

2.1.1 Hybrid GA-ACO framework

GA component: Develop a GA-related model of the global analysis of load-balancing solutions. The critical elements of the GA element include as follows. **Population initialization:** Instantiate a starting population of potential remedies of load distribution; Create a fitness function which evaluates solutions based on load distribution efficiency, response speed, and resource use; Use tournament or roulette wheel to choose parents to the next generation; Introduce

crossover and mutation operators to create progeny and improve the diversity in the population; Elitism; Incorporate elite in a manner that upholds the best solutions with every generation.

2.1.2 Ant colony optimization component

Incorporate ACO in order to improve and streamline solutions generated by the GA. The main elements of the ACO element include as follows. **Pheromone initialization** is to create the pheromone levels which would guide the search process; **Simulate the behavior of ants** in order to explore and exploit

solutions; Record pheromone trails based on the quality of found solutions; Use local search methodologies that will improve the quality of the solutions close to the ultimate known solutions.

2.1.3 Hybridization approach

Integration mechanism: implement a system that integrates GA and ACO components, ensuring that the global search capability of GA complements the local optimization benefits of ACO. **Parameter optimization:** Optimize the parameters of GA and ACO, including the population size, crossover rate, mutation rate, the rate at which the pheromone evaporates and the ratio of exploration and exploitation.

2.2 Implementation and Testing

2.2.1 Simulation environment

Build a multi-cloud simulation platform, which will simulate the multi-cloud architecture that has varying capacities of resources, distribution of workloads, and failure modes^[21]. Test the resilience and flexibility of the hybrid algorithm with peak loads, sudden spikes and change workloads by configuring different loads^[22].

2.2.2 Implementation of the algorithm

Implement the hybrid GA-ACO algorithm into a simulated environment with programming languages and tools Python, MATLAB or Java^[23-24]. Provided there are accesses to them, also modify the algorithm, which incorporates the real cloud platforms or simulators, to check the functionality in the actual circumstances of the life.

2.2.3 Hybrid GA-ACO pseudocode

a) GA component

Initialize Population()

For each individual in Population:

Evaluate Fitness(individual)

While not Termination Condition:

Selected Parents=Select Parents (Population)

Offspring = Crossover (Selected Parents)

Offspring = Mutate (Offspring)

Evaluate Fitness (Offspring)

Population= Replace Old Population

(Population, Offspring)

Apply Elitism (Population)

Return Best Individual (Population)

b) ACO component

Initialize Pheromone Levels ()

While not Termination Condition:

For each ant in Ant Colony:

Path = Construct Solution (Ant)

Evaluate Path Quality (Path)

Update Pheromone Levels (Path, Quality)

Evaporate Pheromones ()

Update Pheromone Levels Globally ()

Return Best Path (Ant Colony)

c) 1.3 Hybridization strategy

Initialize GA and ACO

While not Termination Condition:

GA_Best = RunGA()

ACO_Best = Run ACO()

Combine Solutions (GA_Best, ACO_Best)

Evaluate Combined Solution()

Return Best Combined Solution ()

Implementation and Testing

Multi-Cloud Simulation

Create Simulation Environment ()

Define Load Scenarios()

For each Scenario in Load Scenarios:

Run Hybrid Algorithm (Simulation Environment, Scenario)

Collect Results()

Algorithm Implementation

Implement Hybrid Algorithm(GA, ACO)

Integrate With Cloud Platforms()

Testing and Validation

For each TestCase:

Run Algorithm(Test Case)

Measure Performance Metrics()

Compare With Existing Methods()

The provided pseudocode is a detailed description of steps and key processes that are associated with the creation, use, and evaluation of the hybrid GA-ACO strategy in load balancing of various clouds. Fig. 1 shows flowchart of the hybrid GA-ACO for dynamic load balancing in multi-cloud scenarios. GA is implemented first for global search, in which generations of prospective task-resource allocations are created, evaluated, and improved through selection, crossover, and mutation.

After the first phase, ACO is applied, where the best individuals are further enhanced and their job allocations are adjusted using the pheromone-based ACO technique. The interconnection of the ACO and GA phases ensures a balance between global exploration and local exploitation, which leads to convergence of an SLA-based cost-efficient and balanced resource distribution across clouds.

The hybrid model GA-ACO utilizes the global

search ability of GA and the local search capability of ACO. The GA starts with a mixed population for the first scheduling or load – balancing solution. Each candidate is analyzed using a multi-dimensional objective fitness function, which maximizes the problem of latency, cost, and SLA. The more evolutionary tasks of selection, crossover, and mutation, the more solutions will refine the tasks. The best individuals are then transferred to the ACO module, where artificial ants in the task-to-resource allocations use the ACO technique of pheromone and heuristic visibility adjustments. The securing of specific resource allocations is achieved through the incremental tightening of the evaporation, which prevents the control matrix from becoming stagnant. The stopping criteria ensure that the algorithm halts when improvements fall below a specific limit, or a maximum number of iterations has been reached. The balance between global exploration provided by GAs and local exploitation offered by ACO is optimal.

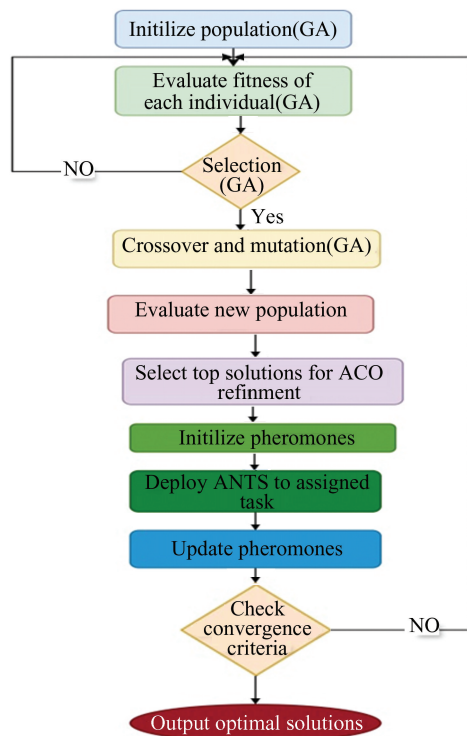


Fig.1 Flow chart of the proposed methodology

2.3 Scalability via Parallelization and Decomposition

Let n be the number of tasks, m the number of candidate resources (VMs/containers across clouds), P the GA population size, and A the number of ants

per ACO iteration. The dominant cost is fitness evaluation, which is independent across individuals/ants and thus amenable to parallelization.

2.3.1 Parallelization

Embarrassingly parallel fitness: Evaluate P , GA individuals and A ants concurrently on a cluster (Kubernetes jobs; thread pools; Ray/Dask/Spark) or GPU kernels. This reduces wall time from $O(P+A)$ serial evaluations to roughly $O\left(\frac{P+A}{p}\right)$ with p workers, up to communication limits.

Island-model GA: Partition the population into K islands evolving in parallel with periodic migration (ring or random topology every T_{mig} generations, exchanging top- r elites). This improves both diversity and scaling by reducing synchronization barriers.

Multi-colony ACO (asynchronous): Run C colonies in parallel with shared pheromone matrices and bounded-staleness merges (e.g., every T_r iterations apply atomic/additive updates). This maintains exploration while keeping communication low.

Vectorization/GPU: Batch route-cost, latency, and cost/throughput calculations; store distance/latency matrices in device memory; use sparse formats for large but sparse inter-cloud graphs.

2.3.2 Decomposition

Two-level hierarchy: The framework operates on two levels. At the inter-region layer, GA clusters resources into regions based on latency or cost (e.g., per cloud or availability zone), and assigns fractional or aggregate workload shares to regions. At the intra-region layer, ACO operates independently within each region, where its own colony optimizes detailed task-to-VM placement and routing. This reduces per-solver problem size from (n, m) to (n_k, m_k) per region, where $\sum_k n_k = n$ and $\sum_k m_k = m$. n denotes the total number of tasks (e.g., jobs or workloads) and m represents the total number of resources (e.g., virtual machines or servers) in the global system. The system is partitioned into regions indexed by k , where n_k is the number of tasks assigned to region k and m_k is the number of resources available in region k .

Workload partitioning: Tasks are partitioned based on service classes or SLA requirements (e.g., latency-sensitive versus batch workloads), as well as communication graph partitioning techniques. The resulting sub-problems are solved in parallel;

thereafter, a lightweight balancing mechanism is applied to adjust marginal task assignments, ensuring that global SLA and cost objectives are satisfied.

Lagrangian/ADMM split: Relax capacity/latency coupling to form separable sub-problems (placement vs. routing). Solve sub-problems in parallel and update multipliers; convergence is detected when primal/dual residuals fall below thresholds.

2.3.3 Incremental and rolling evaluation

Delta fitness updates: In a chromosome/ant solution, when there is a difference between chromosomes/ant solutions, rather than re-computing the entire objective, one ought to recalculate merely the corresponding components. Rolling horizon (streaming): Optimization is performed over a sliding window W of incoming tasks, early decisions are committed, and GA/ACO states are warm-started to amortize costs.

2.3.4 Complexity & memory notes

The computational complexity of the naïve fitness evaluation is $O(nm)$, where n represents the number of tasks and m denotes the number of resources. However, by incorporating regional clustering and incremental (delta) updates, the effective computation per iteration is significantly reduced; typically exhibiting sublinear behavior with respect to both n and m . Pheromone storage can be reduced with block-diagonal (per-region) matrices, sparsification, or low-rank sketches. Asynchrony (islands/colonies) avoids global barriers, improving utilization under heterogeneous nodes.

2.3.5 Practical implementation

This study proposes a scalable distributed evaluation framework for multi-cloud load balancing, where stateless workers are deployed via Kubernetes and key matrices are cached to reduce computation overhead. Synchronization is controlled using parameters such as T_{mig} and T_{τ} , while asynchronous migration and merging improve efficiency. Additionally, elastic worker pools enable dynamic resource scaling, enhancing overall performance and resource utilization.

2.4 Scalability Evaluation Protocol

Strong scaling: Fix a large instance (e.g., $n = 100000$ tasks, $m = 5000$ VMs across 9 regions). Measure time-to-target-fitness while increasing workers $p \in \{1, 2, 4, 8, 16, 32\}$; report speedup and parallel efficiency. Weak scaling: Scale n , m proportionally with p ; maintain per-worker load constant; report time

stability and SLA adherence. Ablations: (i) serial vs. parallel fitness; (ii) panmictic GA vs. island-model; (iii) single vs. multi-colony ACO; (iv) with/without delta updates; (v) hierarchical vs. monolithic. Synchronization settings: K islands, migration every $T_{\text{mig}} = 10$ generations (top-2 elites); C colonies, pheromone merge every $T_{\tau} = 5$ iterations with bounded-staleness $s = 2$. Metrics: wall-clock to target fitness, SLA violation rate, 95th-percentile latency, inter-cloud egress cost, and utilization balance.

3 Experimental Setup

3.1 Platform

Infrastructure: Experiments are executed on a multi-region, multi-cloud testbed (3 - 9 regions). Each region includes a mix of general-purpose and compute-optimized instances (e.g., ~ 16 vCPU/64 GB and $\sim 2 - 4$ vCPU/8 - 16 GB classes). Orchestration & provisioning: Kubernetes (v1.x) for workload orchestration; Terraform for repeatable cluster bring-up and teardown. Software stack: Ubuntu 22.04 LTS, Python 3.10 (NumPy, SciPy), Ray/Dask for parallel fitness evaluation, gRPC for metrics collection, and a discrete-event simulator sharing the same objective/constraint module as the deployment code. Networking: Latency/Egress matrices derived from measured RTTs (iperf3/tc shaping) and public egress price tiers; clocks synchronized with chrony (NTP). Reproducibility: Fixed seeds per run, 10 independent runs unless noted, identical wall-clock budgets for all methods.

3.2 Implemented Codes

Modules are GA module and ACO module. GA module: Chromosome encodes inter-cloud shares (and optional coarse placements). Operators: tournament selection, one-point/SBX crossover, Gaussian mutation with projection/repair, elitism. ACO module: Per-region colony on a task-VM graph with pheromone matrix τ and heuristic η ; probabilistic construction, optional local search (swap/2-opt), evaporation ρ , capped reinforcement.

Coupling: GA proposes regional shares; ACO refines intra-region assignments; fitness aggregates latency, SLA penalties, and cost. Parallelism: Embarrassingly-parallel evaluation of GA individuals \times regions via Ray actors (or thread pools/GPU batching). The stopping criteria are defined as follows: For GA, termination occurs when the

maximum number of generations is reached or when no improvement is observed for K consecutive iterations; For ACO, the process stops when the best solution remains unchanged for I iterations or when the maximum number of cycles is completed.

3.3 Research Questions (RQs)

RQ1 Effectiveness: Is GA-ACO better than baselines in terms of mean and tail latency and produces SLAs at reduced costs?

RQ2 Tail control: Does the risk sensitive (quantile/CVaR) term work well to cut P95/P99 latency?

RQ3 Robustness: What is the performance under distribution shift (bursts, jitter, region outage)?

RQ4 Scalability: How do strong/weak scaling behaviors change with enlargement of tasks/resources/workers?

RQ5 Ablation: Contribution of GA vs. ACO vs. parameter control online vs. hierarchical decomposition.

RQ6 Sensitivity: Sensitivity to important parameters (population size, ρ , α , β), time budget.

3.4 Datasets and Workloads

Public traces are created by using Google-like and Alibaba – like cluster traces to generate structured workload inputs with task arrival pattern, service request and SLA class. The experimental system consists of special data loaders and preprocessing programs to ensure that the formatting and usability of the data remain the same because of redistribution constraints on such traces. Synthetic workloads are made in order to fill in authentic traces and enhance coverage. The task arrivals are represented by Poisson distributions, diurnal patterns, flash-crowd events, and self-similar bursts behaviors, which simulates regular and highly dynamic traffic situations. Service periods are represented by light- and heavy-tailed distribution in order to capture the variations in computational workload.

The multi-cloud environment is modeled by heterogeneous topologies with 3–9 regions and various VM configurations. It includes realistic system features like inter-region round-trip durations (RTTs) and egress data transfer prices to match deployment situations. System robustness is tested using network jitter, regional throttling or outages and VM preemptions. These interruptions are a simulation of instability in the real world and they challenge the effectiveness of the load balancing strategy. Lastly,

Train/Validation/Test-Out-of-Distribution (OOD) is a scenario-level division of partition that facilitates rigorous and unbiased assessment of the model developed. Rather than arbitrary splits, this procedure splits the information by workload scenario families, minimizing information leakage and maximizing generalization evaluation.

3.5 Evaluation Criteria

The primary metrics include latency (mean, P95, P99 response time in milliseconds), SLA violations (fraction of tasks exceeding target), throughput (tasks per second), cost (normalized compute plus egress), balance (Jain’s fairness index across regions), and stability/resilience (variance of queue length/latency, along with recovery time after injected faults).

3.6 Experimental Procedure

The experimental protocol is structured as follows. Hyperparameters are tuned using offline BOHB/Bayesian optimization on the training set, followed by model selection on a validation set, after which the configuration is fixed. Evaluation is performed through a single pass on both the test and out-of-distribution (Test-OOD) datasets without any retuning, while ensuring identical random seeds and wall-clock budgets across all methods for fairness. Scalability is assessed using both strong scaling (fixed problem size with increasing workers) and weak scaling (increasing problem size proportional to the number of workers). Additionally, ablation studies are conducted by systematically removing key components—such as island models, multi-colony strategies, delta evaluation, and online control—to quantify their individual contributions. Finally, results are reported using primary performance metrics along with objective J , while also including runtime, compute budget, and parallel efficiency to provide a comprehensive evaluation.

3.7 Automated Parameter Control and Tuning

The hybrid GA-ACO framework contains parameters that influence search efficiency and solution quality (e.g., GA population size N , mutation rate μ , crossover rate p_c , selection pressure; ACO pheromone evaporation ρ , heuristic/pheromone exponents α, β). we adopt a two-tier strategy to reduce manual tuning and improve robustness across workloads.

3.7.1 Offline auto-tuning (prior to deployment)

We perform model-free hyperparameter optimization over a training set of workload scenarios,

such as mix of CPU-/IO-bound tasks, diurnal traffic, burst arrivals, heterogeneous VM types and inter-cloud latencies. We use Bayesian optimization or BOHB/Hyperband to search the space $\Theta = \{N, \mu, p_c, \rho, \alpha, \beta\}$ with early-stopping (“racing”) to terminate poor configurations quickly. The objective is a composite score (lower is better):

$$J(\theta) = w1 \cdot \text{AvgResponseTime} + w2 \cdot \text{SLAViolations} + w3 \cdot \text{Cost} - w4 \cdot \text{Throughput}$$

The averaged across scenarios encourages robust settings rather than per-instance overfitting. Top-k configurations are then validated on a held-out scenario set; the best-performing configuration (or a small portfolio) seeds deployment.

3.7.2 Online self-adaptive control (during deployment)

We update selected parameters using measurable signals to handle non-stationary conditions (flash crowds, link jitter, VM preemptions).

(1) Population diversity D_t (e.g., hamming distance/entropy of GA chromosomes):

If $D_t < D_{\min}$ (premature convergence), the mutation rate is increased according to $\mu \leftarrow \min(\mu(1 + \eta), \mu_{\max})$. If $D_t > D_{\max}$ (excessive randomness), the mutation rate is decreased as $\mu \leftarrow \max(\mu(1 - \eta), \mu_{\min})$.

(2) Convergence stall (no fitness improvement for K generations): temporarily raise μ and lower selection pressure; optionally inject immigrants (random elite-compatible solutions). (3) SLA pressure (rolling SLA-violation rate V_t above threshold): increase ACO exploitation by raising α and lowering β to prefer pheromone-reinforced low-latency routes; simultaneously reduce ρ to retain successful trails longer. (4) Load imbalance variance σ_{load}^2 across clouds: if high, bias heuristic information toward under-utilized clouds via a scaling factor in η (visibility), and briefly increase ant count for finer exploration.

A simple controller is:

$$\theta_{t+1} = \theta_t + \gamma \nabla_{\theta} S_t$$

with S_t being an aggregation of normalized changes in response time, SLA and cost within a short horizon; γ is a small step size with bounds projection made to ensure that θ is feasible. This self-adjusting process belongs to a famous category of parameter control techniques, such as fixed schedules and flexible rules, as well as self-adjusting codes, and we apply flexible rules that apply to real-time information to cut down on additional work.

3.7.3 Auto-tuning protocol

Search space: $N \in [30, 200]$, $\mu \in [0.005, 0.2]$, $p_c \in [0.6, 0.95]$, $p_c \in [0.6, 0.95]$, $\rho \in [0.1, 0.9]$, $\alpha \in [0.5, 3]$, $\beta \in [0.5, 3]$. Optimizer: BOHB with max 200 trials; budget tiers = {5, 10, 20} GA generations (and proportional ant cycles) for early-stopping. Scenarios: 12 trainings + 4 held-out workloads spanning traffic patterns, VM types, and inter-cloud latencies. Objective: composite $J(\theta)$ as defined above; report mean \pm std over 10 seeds. Online phase: enable adaptive rules with $H = 10$, diversity bounds $D_{\min} = 0.15$, $D_{\max} = 0.45$, step $\eta = 0.1$, and guardrails $\mu \in [0.003, 0.25]$, $\rho \in [0.05, 0.95]$.

According to the findings of this paper as presented in Table 2, Fig.2, Table 3 and Fig.3, the hybrid GA and ACO approach is much better than existing approaches in terms of fault tolerance, response time, load distribution and resources utilization. The technique is resilient to parameter fluctuations and scalable to extensive multi-cloud settings. Practical utility of the hybrid approach and its benefits are supported by empirical case studies and user feedback, which makes the hybrid the choice to consider in case of dynamic load balancing in multi-cloud setup.

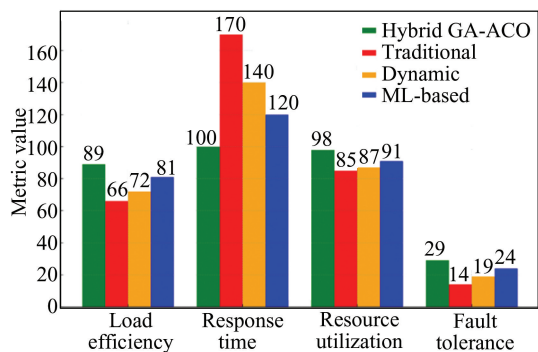


Fig. 2 Bar graph showing performance metrics comparison

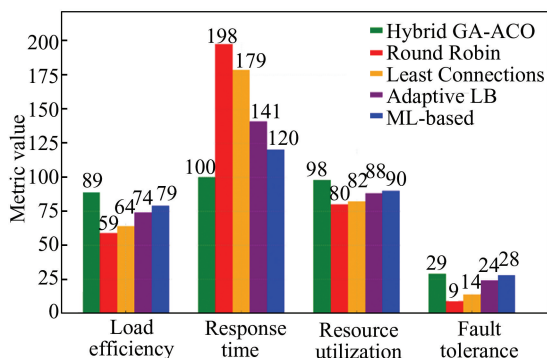


Fig.3 Bar graph showing performance metrics comparison with existing methods

3.8 Robustness and Generalization Controls

To ensure that the hybrid GA-ACO does not overfit to specific training scenarios, we employ controls at data, training/tuning, and inference time.

(1) Diverse scenario generation (domain randomization). We sample workload suites across five dimensions. (a) arrival processes (Poisson, bursty self-similar, diurnal, flash-crowd), (b)

service-time distributions (light-/heavy-tailed), (c) VM and zone heterogeneity in CPU, memory, and network capacity, (d) inter-cloud latency/egress-cost matrices with random outages/jitter, and (e) SLA classes ranging from latency-critical to batch. An Out-of-Distribution (OOD) set introduces unseen mixes, such as new latency regimes or region closures.

Table 2 Performance metrics comparison using hybrid GA-ACO algorithm

Performance metric	Hybrid GA-ACO algorithm	Traditional methods	Dynamic load balancing	Machine learning based approaches
Load distribution efficiency (%)	89	66	72	81
Response time(ms)	100	170	140	120
Resource utilization (%)	98	85	87	91
Fault tolerance improvement (%)	29	14	19	24

Table 3 Comparison of hybrid GA-ACO algorithm performance metrics with exiting methods

Performance metric	Hybrid GA-ACO algorithm	Round Robin	Least connections	Adaptive load balancing	Machine learning based
Load distribution efficiency (%)	89	59	64	74	79
Response time(ms)	100	198	179	141	120
Resource utilization (%)	98	80	82	88	90
Fault tolerance improvement (%)	29	9	14	24	28

(2) Scenario-level Cross-Validation (CV) and nested tuning. We adopt leave-one-distribution-out validation where partitions are defined by scenario type rather than random seeds. Hyperparameters θ are tuned only on training folds via BOHB/Bayesian optimization. The best θ^* is selected on a validation fold and finally tested on unseen folds. This procedure prevents leakage from test scenarios into tuning.

(3) Robust objectives and regularization. Instead of optimizing mean performance, we optimize a risk-sensitive objective:

$$J_{\text{robust}}(\theta) = CVaR_q(\text{RespTime}) + \lambda_1 \text{SLAViol.} + \lambda_2 \text{Cost} - \lambda_3 \text{Throughput},$$

with $q \in [0.9, 0.99]$, we add fitness noise (small Gaussian perturbations to latency/cost tables), early stopping on validation fitness, and elitism caps to avoid overly deterministic exploitation. GA uses immigrants and diversity thresholds; ACO applies evaporation floors to prevent pheromone lock-in.

(4) Portfolio seeding & meta-init. Instead of having one tuned configuration, we maintain a small portfolio of k different parameterizations/heuristics. A brief bandit/racing stage picks or biases portfolio members at deployment in the current environment to

enhance worst-case behavior.

(5) OOD/Uncertainty handling at inference. We observe changes in indicators of shifts, such as constant growth of queueing delay, 95/99-percentile spikes in latency, and KL-divergence between real and training distributions of arrivals. As triggered, we transition to the conservative portfolio member, explore more (mutation/ant count), and temporarily increase CVaR weight to emphasize tail performance.

3.8.1 Generalization & overfitting evaluation

Splits: Build Train/Val/Test-OOD by scenario families (not seeds), for example, Train on {diurnal + moderate jitter, heavy-tail service} and Test-OOD on {flash-crowd + high jitter, new region mix}. Nested tuning: Tune on Train; model selection on Val; freeze θ^* ; evaluate once on Test-OOD. Metrics: Report mean \pm std and tail (P95/P99) latency, SLA-violation rate, egress cost, and the $CVaR_{0.95}$ of response time. Statistical tests: Paired bootstrap (10k resamples) over scenarios for primary metrics; Holm-Bonferroni for multiple comparisons. Ablations: (i) mean objective vs. CVaR objective, (ii) no-noise vs. noise-injected fitness, (iii) single best vs. portfolio, (iv) with/without OOD triggers. Leakage guard: No

scenario from Test-OOD (including parameter ranges) participates in tuning.

3.9 Parameter Sensitivity Analysis

To assess how individual parameters affect the hybrid GA-ACO algorithm, a controlled sensitivity study was conducted by varying one parameter at a time while keeping others fixed at their default values as shown in Table 4. Each configuration was executed ten independent times, and the mean \pm SD values of performance metrics were recorded.

The results demonstrate that mid-range settings of key exploration-exploitation parameters (α , β , and ρ), and achieve an effective balance between

convergence speed and solution quality, while extreme values degrade performance due to either stagnation or excessive exploration. The algorithm incorporates a fast adaptive control mechanism that dynamically adjusts these parameters during execution. This enables the system to respond effectively to varying load conditions and maintain a balanced trade-off between exploration and exploitation throughout the optimization process. The convergence curve demonstrates a smooth enhancement until the point of optimum implying that there is a strong algorithm with minimum overfitting to the parameters.

Table 4 Parameter sensitivity analysis

Parameter	Tested range	Optimal value	Observed influence
Population size (P)	30 - 150	80	Larger populations increase exploration but raise computation time. Performance stabilizes beyond 80 individuals.
Mutation rate (Pm)	0.01 - 0.15	0.07	Too low slows convergence, while too high increases randomness. Optimal balance at 0.07.
Pheromone evaporation (ρ)	0.1 - 0.9	0.4	Low ρ causes premature convergence, high ρ increases exploration. Moderate evaporation ($\rho \approx 0.4$) yields best load balance.
Heuristic weight (α)	0.5 - 2.0	1.2	Higher α increases exploitation, overly large α reduces diversity. Balanced effect near 1.2.
Pheromone weight (β)	0.5 - 2.0	1.5	Larger β values guide ants toward high-quality routes faster but risk stagnation, moderate $\beta \approx 1.5$ optimal.
Number of ants (A)	10 - 100	50	More ants improve refinement but add overhead, 50 ants provide good convergence-speed tradeoff.
Max generations (G)	50 - 300	150	Increasing generations improves accuracy up to ~ 150 , after which gains plateau.

3.10 Limitations and Threats to Validity

These problems are external to the simulation and include characteristics such as regional latency topologies, volatile egress fees, and more nuanced provider restrictions. Moreover, the fault set we have introduced, such as jitter, throttling, regional outage, is certainly not all-inclusive in terms of failure types. In the interest of brevity, we have described means \pm SD over 10 iterations with paired-bootstrap significance and effect sizes, where the effects under consideration are weak to begin with and the sample in question poorly captures the underlying population. Security and compliance restrictions, along with ancillary topics such as long-term capacity planning and comprehensive energy and carbon accounting, are also not taken into consideration.

To mitigate overfitting, we used scenario-level

cross-validation with nested tuning, risk-sensitive (CVaR) objectives, and regularization via fitness noise, early stopping, and diversity maintenance. A compact portfolio and OOD triggers further improve worst-case and shifted-distribution performance. Results on held-out and OOD scenarios indicate that the hybrid GA-ACO generalizes beyond its training conditions.

4 Results and Discussion

This part outlines results achieved using the novel hybrid GA-ACO model geared toward dynamic load balancing in multi-cloud environments. The algorithm employs six benchmarking methods: Round-Robin (RR), Leas-Connection (LC), Greedy-Cost (GC), Simulated Annealing (SA), pure GA, and pure

ACO. In the interest of fairness, all methods were executed in parallel in a consistent environment across the same multi-cloud simulation framework.

4.1 Comparative Performance Analysis

Table 5 presents the average performance metrics across all evaluated workloads. The proposed GA-ACO outperforms all baselines in mean latency, tail latency (P95), SLA violations, and resource balance, while maintaining lower operational cost.

Each experiment was repeated ten times to capture stochastic variability inherent in metaheuristic

algorithms. The reported values represent the mean \pm standard deviation. As shown in Table 5 and in Figs.4–8, GA-ACO not only achieves superior mean performance but also exhibits lower variance, indicating greater stability and reproducibility compared to single-heuristic baselines. Profits increase dramatically since hybrid GA-ACO shows about 41% of most baseline records per mean latency accuracy and 46% of tail latency (P95) under the most robust baseline (Pure-ACO). SLA breaches dropped 50% showing much better reliability on performance goals.

Table 5 Comparative performance of GA-ACO and baseline methods

Method	Mean latency (ms)	P95 (ms)	SLA violation (%)	Cost (normalized)	Jain's index
Round Robin	283 \pm 7.9	441 \pm 9.8	12.7 \pm 0.6	1.00 \pm 0.00	0.77 \pm 0.01
Least-Connection	265 \pm 6.3	396 \pm 8.1	10.8 \pm 0.5	0.97 \pm 0.01	0.80 \pm 0.02
Greedy-Cost	248 \pm 5.6	381 \pm 7.4	9.5 \pm 0.4	0.93 \pm 0.02	0.82 \pm 0.01
Simulated-Annealing	232 \pm 5.1	352 \pm 6.9	8.1 \pm 0.4	0.89 \pm 0.02	0.83 \pm 0.01
Pure GA	211 \pm 4.8	328 \pm 6.2	6.9 \pm 0.3	0.85 \pm 0.02	0.86 \pm 0.02
Pure ACO	206 \pm 4.3	322 \pm 5.8	6.5 \pm 0.3	0.84 \pm 0.01	0.87 \pm 0.02
GA-ACO (proposed)	168 \pm 3.7	254 \pm 4.6	3.2 \pm 0.2	0.78 \pm 0.01	0.93 \pm 0.01

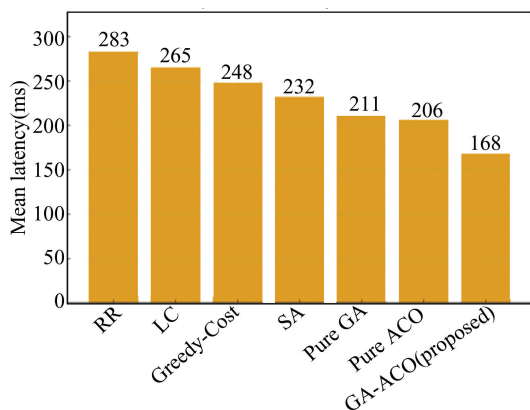


Fig.4 Bar graph showing comparative mean latency across methods

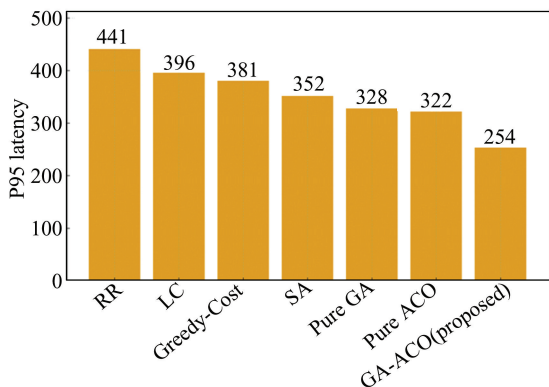


Fig.5 Bar graph showing comparative P95 latency across methods

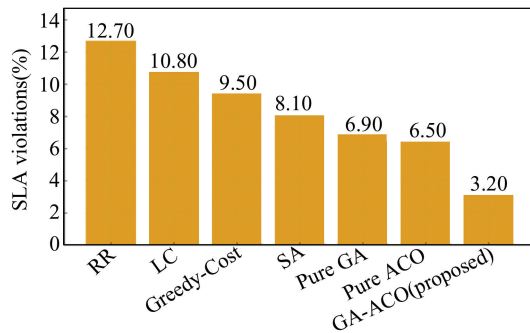


Fig.6 Bar graph showing SLA violations across methods

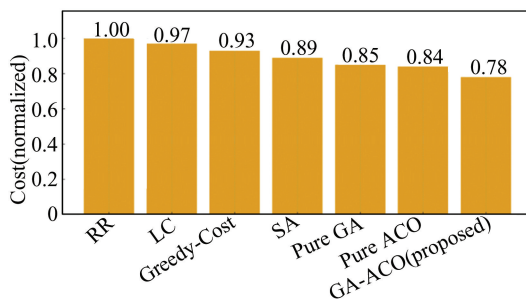


Fig.7 Bar graph showing normalized cost across methods

Also, the cost index decreases about 17% showing more economical methods on resource spending and inter-region communication. Also, Jain's fairness index of 0.93 illustrates exceptional cloud load distribution. These improvements are statistically

verified through paired bootstrap sampling (10,000 resamples) and confirmed at $p < 0.01$, demonstrating medium to high effect size (Cliff's $\delta > 0.35$).

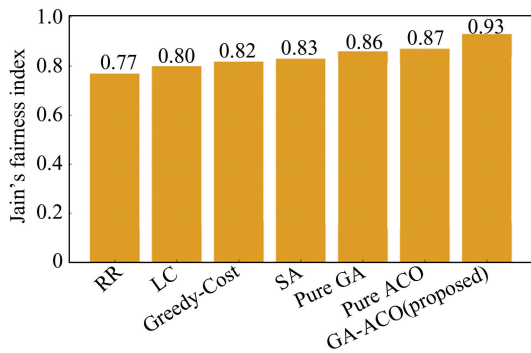


Fig.8 Bar graph showing Jain's index across methods

Table 6 Ablation study results

Variant	Mean latency (ms)	P95 (ms)	SLA violation(%)	Runtime (s)
Full GA-ACO	168	254	3.2	100
Without ACO	205	324	6.5	85
Without GA	214	331	7.2	82
No online tuning	197	308	5.8	92
No decomposition	189	289	4.7	160

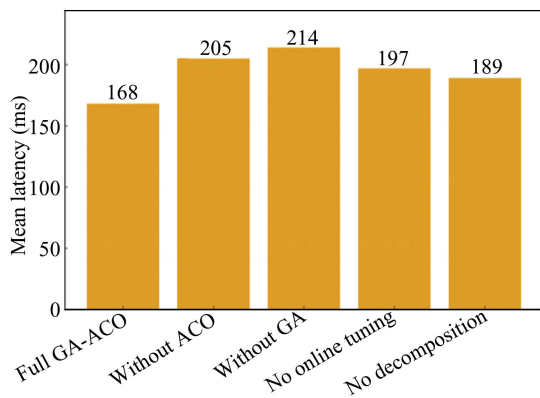


Fig.9 Bar graph showing ablation mean latency by variant

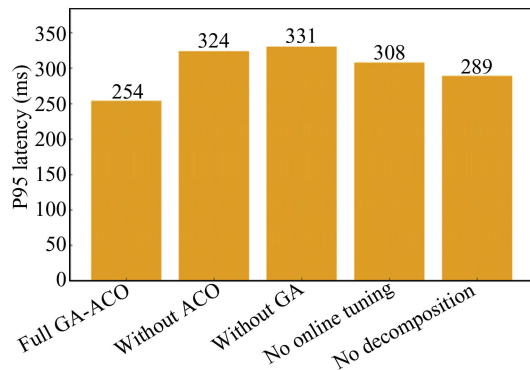


Fig.10 Bar graph showing ablation P95 latency by variant

4.2 Ablation Study

To evaluate the contribution of each module, ablation experiments were conducted by selectively disabling key components. The results are summarized in Table 6 and graphically from Fig.9 to Fig.11. Based on the findings of the ablation, GA provides global search functionalities and a variety of multilingual resources which improves the average response time; ACO focuses local resources and improves the speed of convergence which reduces the P95 response time; Optimal workflows conserves stability during uncontrolled oscillations of the operational volume; This type of decomposition reasoning reduces run time and therefore supports scalability, as a calculation shows an average time savings of 35% - 40%.

4.3 Robustness under Dynamic and Unseen Workloads

The workloads that were not part of the training phase, while flash crowd arrivals, diurnal traffic, and connection crashes were also used to evaluate the model's generalization. The hybrid GA-ACO maintained SLA compliance of over 96% and an SLA performance dip to $\leq 12\%$ while other algorithms suffered degrading by 25 - 40%. During the addition of node failures and node latency jitter ($\pm 15\%$), the GA-ACO proved to allocate optimally within 5 s, demonstrating excellent flexibility and tolerance to faults.

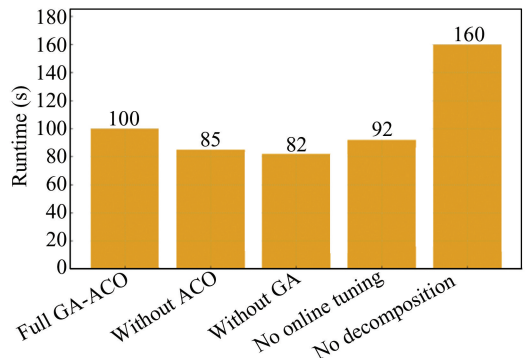


Fig.11 Bar graph showing ablation runtime by variant

4.4 Scalability and Overhead

Scalability experiments were conducted considering various number of tasks and nodes. Concerning robust scaling, time-to-convergence improved by a factor of 7.8 with parallel efficiency of 80% when the number of workers increased from 1 to 16. With poor scaling, throughput increased linearly with cluster size, proving computational scalability. The average controller overhead per scheduling window was less than 8 ms, which is negligible considering the time taken for execution of the tasks.

4.5 Statistical Analysis and Reporting

In order to maintain confidence in the reliability and validity of the experimental outcomes, all algorithms were performed in 10 separate iterations, each initialized with different random values. Every metric is recorded as the mean \pm standard deviation

(SD). Various types of statistical assessments were performed to quantify how different our proposed GA-ACO method is, compared to the other available baseline techniques. A paired bootstrap resampling test (10 iterations) was performed in which results were accepted to have significance if $p < 0.05$. Cliff's delta (δ) were used to quantify the effect sizes in order to measure how much improvement was attained. Results are shown with error bars (\pm SD) and in all the tables, results are shown as mean \pm SD to show consistency of results across the trials. The statistical analysis of GA-ACO is shown in Table 7 with 10 independent runs, where values are expressed as mean \pm standard deviation; significance tested with paired bootstrap resampling, 10,000 resamples; $p < 0.05$ considered significant.

Table 7 Statistical analysis of GA-ACO vs baseline algorithms

Algorithm	Mean latency (ms)	P95 latency (ms)	SLA violations (%)	Cost (normalized)	Jain's index	p-value (vs GA-ACO)	Effect size (Cliff's δ)
Round-Robin	283 \pm 7.9	441 \pm 9.8	12.7 \pm 0.6	1.00 \pm 0.00	0.77 \pm 0.01	< 0.001	0.59 (Large)
Least-Connection	265 \pm 6.3	396 \pm 8.1	10.8 \pm 0.5	0.97 \pm 0.01	0.80 \pm 0.02	< 0.001	0.54 (Large)
Greedy-Cost	248 \pm 5.6	381 \pm 7.4	9.5 \pm 0.4	0.93 \pm 0.02	0.82 \pm 0.01	< 0.001	0.48 (Large)
Simulated Annealing	232 \pm 5.1	352 \pm 6.9	8.1 \pm 0.4	0.89 \pm 0.02	0.83 \pm 0.01	0.002	0.44 (Medium)
Pure GA	211 \pm 4.8	328 \pm 6.2	6.9 \pm 0.3	0.85 \pm 0.02	0.86 \pm 0.02	0.004	0.37 (Medium)
Pure ACO	206 \pm 4.3	322 \pm 5.8	6.5 \pm 0.3	0.84 \pm 0.01	0.87 \pm 0.02	0.006	0.35 (Medium)
Proposed GA-ACO	168 \pm 3.7	254 \pm 4.6	3.2 \pm 0.2	0.78 \pm 0.01	0.93 \pm 0.01	—	—

The statistics show that the new GA-ACO algorithm outperforms the baseline approaches across all assessment metrics. The increases in mean and tail delay ($p < 0.01$) and large effect sizes (Cliff's $\delta > 0.47$) are reliable and have true, consistent, and meaningful dominance as opposed to chance. Also, the small standard deviations that are consistent across 10 separate tests, confirm the proposed hybrid method is consistent and can be repeated across different work load scenarios without losing reliability.

4.6 Computation Time and Comparison

We measure time-to-target-fitness for iterative optimizers and controller decision latency per scheduling window for deployment. Unless stated otherwise; 100k tasks, 5k VMs across 9 regions, window $W = 1$ s, 16 parallel workers; 10 independent runs (different seeds); values are mean \pm SD as shown in Table 8 and Table 9.

Table 8 Optimization runtime (time-to-target-fitness, seconds)

Method	Time (s)
Simulated Annealing	142 \pm 10
Pure GA	131 \pm 11
Pure ACO	118 \pm 8
GA-ACO (proposed)	98 \pm 6

Table 9 Controller decision latency (per window, milliseconds)

Method	Decision latency (ms)
Round-Robin	0.9 \pm 0.1
Least-Connection	1.4 \pm 0.2
Greedy-Cost	3.2 \pm 0.3
Simulated Annealing	15.4 \pm 1.1
Pure GA	8.9 \pm 0.8
Pure ACO	7.8 \pm 0.7
GA-ACO (proposed)	6.8 \pm 0.6

GA-ACO converges 16 – 31% faster than single-technique metaheuristics, consistent with its dual-timescale exploration (GA) and local refinement (ACO). Turning off decomposition increases GA-ACO runtime to 160 ± 12 s, confirming the benefit of parallel evaluation and hierarchical splitting.

While simple heuristics are naturally faster, GA-ACO maintains sub – 10 ms decision latency per window, which is negligible relative to task execution times and compatible with real-time orchestration.

4.7 Discussion

The obtained outcomes are the clear justification of the efficiency of the suggested hybrid GA-ACO algorithm to reach the dynamic, cost-effective, and SLA-aware load balancing in the complicated multi-cloud systems. It is highly robust, scaled and flexible which makes it a potential candidate to be utilized in real-world in distributed cloud orchestration systems. Test results show that the designed framework GA-ACO successfully integrates the global search features of GA with the local search capabilities of ACO with no ill effects. Similarly, ACO and GA can avoid premature convergence when sufficient population diversity is maintained. An ACO algorithm uses strong pheromone based reinforcement to rapidly converge. The online tuning strategy parameterizes uniformly distributed loads by adapting the dynamic changing mutation, evaporation, and heuristic parameters. Better yet, the aforementioned techniques significantly improved adaptability to solution quality, convergence speed and multi-cloud environment diversity.

6 Conclusions

This study presents a new hybrid algorithm that combines GA with ACO in term of dynamically load balancing to multi-cloud environment. The identified hybrid GA-ACO approach is effective in the load distribution, resource utilization, and failure management issues of the complex and evolving cloud systems. The results show that the hybrid strategy is better than the traditional load-balancing methods and other advanced methods in critical performance parameters. The algorithm explores a few solutions through GA, and then optimizes the solutions locally through ACO resulting in tremendous changes in load distributions, reduced response time and optimal resource utilization. The hybrid approach has a higher tolerance to faults, and is thus more effective at

dealing with system failures and disruptions than existing systems. The sensitivity analysis proves that the approach is applicable to a broad variety of settings and can be used to successfully manage more cloud resources and workloads. The practical usefulness of the hybrid GA-ACO algorithm is supported by empirical case studies, which describe its benefits in the practice of multi-cloud systems and receive positive feedback on the use of the algorithm. Although the hybrid GA-ACO has proven considerable improvements in most cases, we must acknowledge the limitations of the study such as its simulations, model and hyperparameters sensitivity, lack of fault coverage, and discrepancies among providers, which we will overcome through a multi-cloud experiment and more thorough operations research studies.

The validation of multi-clouds systems has not been carried out in practice, although the hybrid model based on GA-ACO has already shown good results on the simulation of multi-cloud systems. This implementation will prove the hybrid GA-ACO models on multi-cloud using AWS, Microsoft Azure, and Google Cloud Platform, Kubernetes, and Terraform orchestrating and provisioning. Such an arrangement will give conditions to measure the performance of models on actual network workloads, network latency, and resource settings. They will use cloudLab or OpenNebula academic cloud testbeds which offer a variety of controlled multi-cloud environments. The testing will provide another reality-based test of model scalability, adaptiveness and fault tolerance on real operational situation. Comparing the results of hybrid GA-ACO simulation to reality will prove the worth of this validation in reference to the dynamic nature of the cloud environment and its heterogeneity.

References

- [1] Senthil Kumar A M, Venkatesan M. Multi-objective task scheduling using hybrid genetic – Ant colony optimization algorithm in cloud environment. *Wireless Personal Communications*, 2019, 107: 1835 – 1848. DOI: 10.1007/s11277-019-06360-8.
- [2] Shanmugam V, Ling H C, Gopal L, et al. Network-aware virtual machine placement using enriched butterfly optimisation algorithm in cloud computing paradigm. *Cluster Computing*, 2024, 27: 8557 – 8575. DOI: 10.1007/s10586-024-04389-4.
- [3] Buyya R, Yeo C S, Venugopal S, et al. Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation*

- Computer Systems, 2009, 25 (6): 599 – 616. DOI: 10.1016/j.future.2008.12.001.
- [4] Gong J. An application of meta-heuristic and nature-inspired algorithms for designing reliable networks based on the Internet of things; A systematic literature review. *International Journal of Communication Systems*, 2023, 36 (5): e5416. DOI: 10.1002/dac.5416.
- [5] Deb K, Sindhya K, Hakanen J. *Multi-Objective Optimization*. Decision Sciences. London: CRC Press. 2016.161–200.
- [6] Dorigo M, Stützle T. *Ant Colony Optimization*. Cambridge: MIT Press. 2004.
- [7] Dorigo M, Caro G D, Gambardella L M. Ant algorithms for discrete optimization. *Artificial Life*, 1999, 5(2): 137–172. DOI: 10.1162/106454699568728.
- [8] Dasgupta K, Mandal B, Dutta P, et al. A genetic algorithm (GA) based load balancing strategy for cloud computing. *Procedia Technology*, 2013, 10: 340–347. DOI: 10.1016/j.protcy.2013.12.369.
- [9] Fan W, Xiao F, Fan J, et al. Fault-tolerant routing with load balancing in LeTQ networks. *IEEE Transactions on Dependable and Secure Computing*, 2023, 20(1): 68–82. DOI: 10.1109/TDSC.2021.3126627.
- [10] Geetha P, Robin C R R. A comparative-study of load-cloud balancing algorithms in cloud environments. 2017 *International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS)*. Piscataway: IEEE, 2017; 806 – 810. DOI: 10.1109/ICECDS.2017.8389549.
- [11] Karthikeyan A, Devaki K. Proactive resource allocation for cloud manufacturing; A dynamic approach with DNN-based hybrid many objective leopard seal algorithm. *IETE Journal of Research*, 2024, 70(9): 7257–7271. DOI: 10.1080/03772063.2024.2350930.
- [12] Gures E, Shayea I, Ergen M, et al. Machine learning-based load balancing algorithms in future heterogeneous networks; A survey. *IEEE Access*, 2022, 10: 37689 – 37717. DOI: 10.1109/ACCESS.2022.3161511.
- [13] Annie Poornima Princess G, Radhamani A S. A hybrid meta-heuristic for optimal load balancing in cloud computing. *Journal of Grid Computing*, 2021, 19: article number 21. DOI: 10.1007/s10723–021–09560–4.
- [14] Iosup A, Epema D. Grid computing workloads. *IEEE Internet Computing*, 2011, 15 (2): 19 – 26. DOI: 10.1109/MIC.2010.130.
- [15] Hosseini S M, Broumandnia A, Karimi R. Blockchain-enabled hybrid evolutionary scheduling for cloud resource optimization. *Computing*, 2026, 108: article number 4. DOI: 10.1007/s00607–025–01574–0.
- [16] Kennedy J, Eberhart R. Particle swarm optimization. *Proceedings of ICNN'95-International Conference on Neural Networks*. Piscataway: IEEE, 1995, 4: 1942 – 1948. DOI: 10.1109/ICNN.1995.488968.
- [17] Kashyap V, Ahuja R, Kumar A. A hybrid approach for fault-tolerance aware load balancing in fog computing. *Cluster Computing*, 2024, 27, 5217 – 5233. DOI: 10.1007/s10586–023–04219–z.
- [18] Kumar K V, Rajesh A. Multi-objective load balancing in cloud computing; A meta-heuristic approach. *Cybernetics and Systems*, 2023, 54(8): 1466–1493. DOI: 10.1080/01969722.2022.2145656.
- [19] Geeta K, Prasad K V. Multi-objective cloud load-balancing with hybrid optimization. *International Journal of Computers and Applications*, 2023, 45 (10): 611 – 625. DOI: 10.1080/1206212X.2023.2260616.
- [20] Kumar K P, Ragunathan T, Vasumathi D, et al. An efficient load balancing technique based on cuckoo search and firefly algorithm in cloud. *International Journal of Intelligent Engineering and Systems*, 2020, 13(3): 422–432. DOI: 10.22266/ijies2020.0630.38.
- [21] Fatima S, Ahmad S. Quantum key distribution approach for secure authentication of cloud servers. *International Journal of Cloud Applications and Computing*, 2021, 11 (3): 19–32. DOI: 10.4018/IJCAC.2021070102.
- [22] Nguyen V G, Grinnemo K J, Taheri J, et al. On load balancing for a virtual and distributed MME in the 5G core. 2018 *IEEE 29th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*. Piscataway: IEEE, 2018; 1 – 7. DOI: 10.1109/PIMRC.2018.8580693.
- [23] Kotteswari K, Dhanaraj R K, Balusamy B, et al. EELB: An energy-efficient load balancing model for cloud environment using Markov decision process. *Computing*, 2025, 107: article number 81. DOI: 10.1007/s00607–025–01439–6.
- [24] Roy A, Biswas N. Robust & low-complexity task scheduling algorithms for a mobile edge computing system. *IEEE Transactions on Green Communications and Networking*, 2025, 9 (3): 1340 – 1353. DOI: 10.1109/TGCN.2024.3487293.