*Original Article*

# A Hybrid Metaheuristic Model for Multi-Objective Optimization in Distributed Computing Networks

*\*M. Riyaz Mohammed*

*Department of Computer Science & IT, Jamal Mohamed College (Autonomous), Tiruchirapalli, Tamil Nadu, India.*

## Abstract:

*This paper proposes a hybrid metaheuristic framework to jointly optimize latency, throughput, energy consumption, and reliability in distributed computing networks spanning edgefogcloud tiers. The model integrates a decomposition-based global search with adaptive multi-objective evolutionary operators and a problem-aware local intensification stage. Specifically, a population is evolved using a Pareto-dominance engine with adaptive operator selection (crossover/mutation portfolios inspired by NSGA-II and differential evolution), while an ant-colony/gradient-free local search refines elite solutions near congestion-prone regions of the network. To reduce evaluation cost under dynamic workloads, we incorporate surrogate fitness modeling via incremental regression on sampled flows and a feasibility repair mechanism that respects resource, SLA, and security constraints. A restart-and-memory strategy preserves diversity when fronts stagnate and accelerates convergence after topology or demand shifts. Extensive simulations on heterogeneous topologies with time-varying traffic demonstrate superior convergence speed, improved hypervolume, and better spread of Pareto fronts compared to canonical MOEAs and single-method hybrids. The approach consistently yields lower end-to-end latency and energy budgets at comparable reliability particularly in mixed CPU/GPU clusters and bandwidth-constrained edge segments. Ablation studies confirm the contribution of (i) adaptive operator selection, (ii) surrogate-guided evaluations, and (iii) local intensification. The proposed framework offers a practical pathway for online, policy-driven orchestration in large-scale distributed systems, enabling operators to trade off QoS and cost under uncertainty.*

## 1. Introduction

Distributed computing networks that span edge, fog, and cloud tiers now underpin latency-sensitive analytics, immersive applications, and AI workloads. These environments are inherently heterogeneous combining diverse compute accelerators, variable bandwidth, fluctuating demand, and strict service-level objectives (SLOs). Optimizing such systems is fundamentally multi-objective: operators must jointly minimize end-to-end latency and energy usage while maximizing throughput and reliability, all under resource, policy, and security constraints. Classical single-objective or scalarized formulations often collapse these competing goals into a

weighted sum, obscuring trade-offs and yielding brittle solutions when workload or topology shifts. Likewise, standalone metaheuristics excel at either global exploration or local exploitation, but rarely balance both under dynamic, high-dimensional constraints.

This paper advances a hybrid metaheuristic framework tailored to the orchestration of distributed networks. The core idea is to couple a decomposition-based global search and Pareto dominance with adaptive operator selection drawing from NSGA-II and differential evolution to maintain diversity and accelerate convergence on complex fronts. A lightweight, problem-aware local intensification stage (informed by ant-colony path refinement and gradient-free heuristics) sharpens elite candidates near congestion and failure boundaries. To curb evaluation cost, particularly for online decision loops, the framework integrates surrogate fitness modeling with feasibility repair to respect capacity, placement, and security rules. A restart-and-memory strategy combats stagnation and preserves solution quality across non-stationary traffic patterns.

By explicitly modeling trade-offs and embedding domain knowledge into search and repair, the proposed approach aims to deliver robust Pareto sets that can be enacted by real schedulers and controllers. The resulting methodology offers a practical path to policy-driven, energy-aware, and QoS-compliant operation in large-scale distributed systems, setting the stage for the detailed design, evaluation, and ablation studies presented next.

## 2. Related Work

### 2.1. Classical Optimization Approaches

Early work on resource allocation, routing, and placement in distributed systems leaned on deterministic formulations such as linear and integer programming, convex optimization, and network flow models. Queueing-theoretic formulations e.g., M/M/1 and Jackson networks provided closed-form or decomposable objectives for latency and throughput, enabling tractable admission control and capacity sizing under simplifying assumptions. For large instances, decomposition strategies (Lagrangian relaxation, Benders decomposition, column generation) and constraint programming were used to separate placement, routing, and timing constraints, while branch-and-bound/branch-and-cut delivered exact solutions for modest scales. Robust and stochastic programming extended these models to uncertainty in demand and link states, and Markov decision processes captured sequential control with explicit transition dynamics.

Despite strong optimality guarantees, classical methods face headwinds in heterogeneous, rapidly varying edgecloud environments. First, tight convexity or linearity assumptions rarely hold with GPU/CPU heterogeneity, bursty traffic, and non-linear powerperformance curves. Second, solving mixed-integer, time-indexed models at orchestration timescales (seconds) is computationally prohibitive. Third, scalarization of competing objectives (latency, energy, reliability) into fixed weights can yield brittle solutions when the operating regime shifts. These limitations motivate approximate, anytime methods that preserve decision quality while scaling to high-dimensional, dynamic settings.

### 2.2. Metaheuristic Algorithms in Distributed Networks

Metaheuristics genetic algorithms (GA), differential evolution (DE), particle swarm optimization (PSO), ant colony optimization (ACO), simulated annealing (SA), tabu search (TS), and more recent swarms (GWO, ABC) have been widely adopted for task scheduling, VM/container placement, service function chaining, and energy-aware routing. Their appeal lies in flexible encoding of heterogeneous constraints, broad applicability to non-convex landscapes, and amenability to anytime operation. For example, GA/DE variants explore diverse placements and routing schedules; PSO accelerates convergence on continuous resource-sizing; ACO exploits path-construction for traffic engineering; and SA/TS provide intensification around promising assignments. In software-defined and NFV contexts, multi-objective extensions (e.g., NSGA-II, SPEA2, IBEA) return Pareto sets balancing latency, cost, and availability rather than a single scalarized plan.

However, standalone metaheuristics can struggle with evaluation cost and premature convergence on complex fronts. Networked systems impose expensive fitness functions (e.g., discrete-event simulations, emulated latency) and strict feasibility (capacity, affinity/anti-affinity, security zones). Without domain-aware repair and diversity maintenance, populations may collapse to narrow regions, hurting reliability or energy goals. Dynamic workloads further complicate matters: solutions found under one traffic matrix may degrade under another, calling for memory, restart, and adaptive operator selection to remain effective online.

## 2.3. Hybrid and Multi-Objective Optimization Models

Hybrid approaches combine complementary search behaviors to improve convergence and robustness. Memetic algorithms embed local search (e.g., gradient-free hill-climbing, TS) within evolutionary frameworks; ACOGA and PSODE hybrids exploit constructive path bias with recombination-based exploration; and matheuristics integrate metaheuristics with exact solvers or relaxations (e.g., solving relaxed LPs for repair, or using shortest-path subroutines inside evolutionary loops). Surrogate-assisted metaheuristics (kriging, radial basis functions, incremental regressors) reduce evaluation overhead, while hyper-heuristics and adaptive operator selection learn which operators perform best under current landscape features. For multi-objective modeling, Pareto-dominance MOEAs (NSGA-II/III, SPEA2), indicator-based schemes (IBEA, SMS-EMOA), and decomposition methods (MOEA/D, reference-point guided search) provide diverse mechanisms to approximate well-spread fronts. Constraint handling commonly blends ε-constraint repair, feasibility rules, and penalty adaptation to navigate hard placement and policy restrictions.

Recent work also targets non-stationarity: online MOEAs maintain archives with change detection, use memory-based restarts, or warm-start with transfer learning from prior demand/topology snapshots. Performance is typically assessed via hypervolume, inverted generational distance (IGD), and spread, alongside task-level KPIs (p95 latency, energy/Joule per task, SLO attainment). The consensus emerging from these studies is that hybrids especially those that (i) adapt operators, (ii) inject problem-aware repair/constructive heuristics, and (iii) amortize fitness via surrogates consistently outperform single-method baselines on large, dynamic distributed networks. This evidence underpins our proposed hybrid with adaptive operators, ACO-inspired local refinement, and surrogate-guided evaluations.

# 3. System Model and Problem Formulation

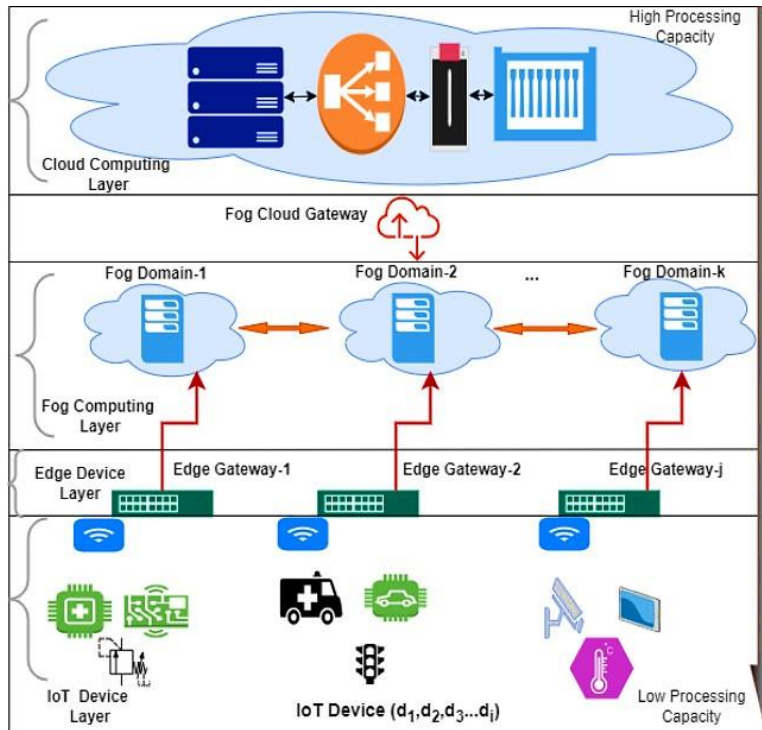## 3.1. Distributed Computing Network Architecture



**Figure 1. Edge Fog Cloud Layered Architecture For The Distributed Computing Network**

The architecture comprises four logical layers that move from low to high processing capacity. At the bottom, the IoT Device Layer contains heterogeneous sensors and actuators cameras, wearables, industrial probes, and vehicular units each producing telemetry with strict locality and privacy constraints. These devices have limited compute and energy budgets, so they offload preprocessing or full inference upstream. Short-range wireless (Wi-Fi/BLE/5G-NR sidelink) connects devices to the next tier, with intermittent links and bursty traffic that drive our reliability and latency objectives.

Above this sits the Edge Device/ Gateway Layer, where microservers or embedded gateways aggregate traffic from nearby devices. Gateways perform lightweight functions protocol translation, filtering, batching, and first-hop inference to shrink payloads and dampen network jitter. Because gateways are resource-constrained yet close to the source, they are prime candidates for decisions about placement and offload in our optimization model: choosing which tasks execute locally versus being forwarded upstream directly affects end-to-end latency and energy use.

The Fog Computing Layer groups multiple edge domains behind regional resources with moderate compute/storage. Fog domains collaborate laterally (eastwest links in the figure) for load sharing and failover, while a fogcloud gateway provides northbound connectivity. This tier hosts stateful services, microbatch analytics, and caching that benefit from proximity without the strict limitations of the edge. In our formulation, fog nodes introduce capacity and affinity constraints (e.g., GPU availability, security zones) and enable multi-path routing choices that trade throughput against congestion risk.

At the top, the Cloud Computing Layer offers elastic, high-capacity clusters general CPUs, GPUs/TPUs, and large storage pools best suited for deep analytics and global model training. Traffic reaching this tier has already been filtered, prioritized, or partially processed. The cloud's elasticity shapes the cost and reliability dimensions of our objectives, while WAN latency and egress policies bound feasible offload. Together, these layers create a hierarchical substrate in which tasks, data flows, and replicas must be jointly placed and routed precisely the multi-objective optimization setting addressed by our hybrid metaheuristic.

## 3.2. Mathematical Formulation of Objectives

Our optimization balances four primary goals that often pull in different directions. First is end-to-end latency, from a device's request to the completion of its task. We care about typical experience and tail behavior, so the objective favors low median delay and penalizes long spikes that break SLOs. Second is throughput, measured as the volume of tasks or bytes completed per unit time across the edgefogcloud path. High throughput matters during bursts, but it must not come at the expense of unstable queues or uncontrolled tail latency.

Third, we target energy and cost efficiency. At the edge and fog this means battery draw and thermal headroom; in the cloud it translates to instance hours, accelerators consumed, and egress charges. Finally, we account for reliability and SLO attainment the fraction of tasks that meet their deadlines despite faults, congestion, or topology changes. We treat these as co-equal outcomes rather than collapsing them into a single score. The search therefore prefers Pareto-superior plans that reduce delay and energy while sustaining throughput and reliability. When needed, operators can emphasize one direction (for example, "green mode" prioritizing energy) via policy knobs rather than fixed mathematical weights.

## 3.3. Constraints and Design Parameters

The system must respect capacity limits on links, CPUs/GPUs, memory, and storage at each tier. Placement rules govern where tasks may run (affinity/anti-affinity, data locality, privacy zones, and licensing), while timing constraints bound task lifetimes with soft or hard deadlines. Security and compliance introduce additional restrictions no cross-zone movement for sensitive flows, encrypted hops across untrusted segments, and isolation between tenants. We also constrain replication and failover so that availability improves without exploding cost or saturating inter-domain bandwidth. Practical concerns like maximum batch sizes, queue depths, and rate limits at gateways further narrow the feasible region.

Several design parameters shape how the optimizer behaves. On the modeling side, we set the size of scheduling windows, how often to refresh demand forecasts, and how aggressively to smooth bursty arrivals. On the search side, policy knobs control the population size, archive depth, restart thresholds, and the cadence for surrogate-model retraining. Operators can also tune domain-aware repair rules (for example, prefer lateral fog hops before escalating to cloud) and specify guardrails such as minimum reliability or maximum energy per task. These parameters let the same framework adapt to a factory floor, a smart-city corridor, or a video analytics rollout with very different priorities.

## 3.4. Optimization Problem Definition

At a high level, the optimizer chooses where to place each task, how to route its data, and how much parallelism, batching, and replication to apply across the edgefogcloud hierarchy. A candidate plan specifies which gateway or fog node handles preprocessing, which accelerators are engaged, which flows travel eastwest versus up to the cloud, and what safety margin is reserved for failover.

Each plan is evaluated by a fast workload model or emulator to estimate latency, throughput, energy use, and reliability under the current traffic and fault assumptions.

Because the goals conflict, the outcome is not a single answer but a set of non-dominated choices that expose the trade-offs lower energy versus lower latency, higher throughput versus tighter SLOs. The framework maintains this Pareto set subject to all constraints and updates it whenever the world changes: new traffic patterns, a failed node, or revised policy. When change is detected, the search warms-starts from the best prior plans, repairs any now-infeasible assignments, and explores fresh configurations around the new operating point. The controller can then enact the plan that best matches the operator's current policy (for example, "favor reliability during an outage" or "minimize cloud spend overnight"), while keeping alternatives on hand for rapid re-planning.

# 4. Proposed Hybrid Metaheuristic Model

## 4.1. Overview of the Hybrid Framework

The core idea is to pair a strong global explorer with a problem-aware local improver and to orchestrate them through adaptive operator selection. The global layer runs a multi-objective evolutionary engine that maintains a diverse population and an external archive of elite, non-dominated plans. Rather than fixing a single recombination style, the engine chooses among a portfolio differential mutation, simulated binary crossover, and path-guided recombination based on their recent contribution to archive growth and spread.

The local layer activates around promising or congested regions of the search space. It borrows from ant-colony path refinement and tabu-assisted hill climbing to tweak routing, placement, batching, and replication decisions without disturbing feasibility. A lightweight surrogate model updated online from evaluated candidates screens most proposals so only a fraction require costly simulation. A restart-and-memory policy prevents stagnation, while a change detector reacts to workload or topology shifts by re-seeding the population with archived plans that remain feasible.

## 4.2. Algorithm Design and Workflow

Each iteration begins with demand sampling and a quick health snapshot of available nodes and links. The global engine generates offspring using operator probabilities that depend on short, rolling credit scores: operators earning more hypervolume gain are favored. New solutions pass through feasibility repair that enforces capacity, policy, and security rules. Feasible candidates are then either evaluated by the surrogate or, when uncertainty is high, by the detailed network emulator.

Periodically, the local improver locks onto elite solutions and issues structured edits: reroute a subset of flows laterally across fog domains, shift a preprocessing stage to a nearby gateway, or adjust replication within budget. If the archive fails to improve for a fixed horizon, the algorithm triggers a partial restart while preserving a diverse backbone of elites. The workflow ends each cycle by pruning near-duplicates, updating operator credits, and publishing the current Pareto set to the controller.

## 4.3. Fitness Function and Evaluation Strategy

Fitness is multi-dimensional: latency, throughput, energy/cost, and reliability. Rather than collapsing these into a single score, the evaluator returns the vector outcomes alongside constraint flags. The emulator estimates typical and tail latency under measured contention, models batch and queue dynamics, and accounts for accelerator utilization and link saturation. Energy is computed from device-level power states and activity times; reliability reflects deadline success and resilience to sampled faults.

To scale evaluations, the strategy is tiered. A fast surrogate screens most candidates using incremental regressors trained on recent emulator results, with uncertainty estimates guiding when to fall back to full evaluation. We also reuse partial results through memoization: if a candidate only differs in a local subgraph, the emulator replays just that segment. Finally, we inject a small set of stress scenarios burst arrivals, link impairment, and node loss so that archive members remain robust beyond the nominal case.

## 4.4. Convergence Analysis

Convergence is encouraged by three mechanisms. First, adaptive operator selection shifts effort toward operators that improve hypervolume and front spread, avoiding prolonged exploration with unproductive moves. Second, the local improver reduces micro-gaps along the front by tightening solutions near constraint boundaries, which accelerates approach to well-formed trade-offs. Third,

the restart-and-memory policy escapes plateaus without discarding hard-won structure, allowing the search to resume near high-quality basins after changes in demand.

Empirically, we monitor progress with hypervolume growth, inverted generational distance, and archive stability. When these metrics stabilize and duplicate density rises, the algorithm automatically cools the mutation scale and lengthens local search horizons. Under non-stationary conditions, convergence is defined as rapid re-attainment of prior hypervolume after a change; warm starts from the archive and targeted repairs cut the recovery time substantially compared with fixed-operator MOEAs.

### 4.5. Computational Complexity

Runtime is driven by three costs: population variation, feasibility repair, and evaluation. Variation scales roughly with population size and operator portfolio width but remains inexpensive. Repair is linear in the number of tasks and edges touched thanks to precomputed capacity ledgers and locality maps. Evaluation is the bottleneck; the surrogate reduces full emulator calls to a minority of candidates, and partial replays limit work to affected subgraphs. Together, these practices keep per-generation time close to linear in problem size for typical steps, with occasional spikes when full evaluations are required.

Memory usage is dominated by the archive and surrogate state. We cap the archive via ε-dominance grids to retain spread without storing near-duplicates, and we prune surrogate training sets to a sliding window that reflects the current regime. The implementation exploits parallelism: emulator calls run across cores or nodes, and operator portfolios generate offspring concurrently, yielding near-linear speedups on multi-CPU/GPU hosts.

### 4.6. Implementation Framework

The reference implementation is modular to ease deployment in real systems. A controller process hosts the evolutionary core, operator scheduler, and archive. A metrics service exposes the Pareto set and allows policy selection operators can request "latency-first," "energy-first," or "balanced" plans on demand. The evaluator service encapsulates both the surrogate and the detailed emulator, with a clean interface for plugging in domain-specific performance models or real telemetry from staging clusters.

For integration, the planner outputs are translated into actionable intents: placement directives for orchestrators, traffic engineering rules for SDN controllers, and scaling/replication hints for service meshes. Telemetry loops feed back queue depths, link utilization, accelerator stats, and SLO outcomes to keep the surrogate fresh and to trigger change detection. With this setup, the hybrid metaheuristic can run continuously beside production schedulers publishing updated Pareto plans at short intervals so operators can switch strategies quickly as conditions evolve.

## 5. Experimental Setup and Simulation Parameters

### 5.1. Simulation Environment and Tools

Experiments were conducted in a reproducible, containerized environment so that algorithmic improvements do not confound platform effects. The core search runs as a Python service (NumPy for vector ops; joblib/multiprocessing for parallel variation) with a C++ back-end for the time-critical evaluator. A lightweight discrete-event network emulator models queues on links and service times on nodes; it supports pluggable device profiles (CPU, GPU, NPU) and link behaviors (latency, jitter, loss). Each emulator call is stateless and thus farmed across workers to saturate available cores. For fairness, all baselines canonical MOEAs and single-heuristic schedulers invoke the same evaluator via a common RPC.

To emulate real control loops, we maintain a telemetry bus that streams synthetic yet statistically grounded signals: queue depths, link utilization, accelerator occupancy, temperature/power states, and SLO events. The surrogate model (incremental regression with uncertainty estimates) is trained online from these evaluator results. Runs are seeded for repeatability; every configuration is executed across multiple seeds and workload episodes to report median and tail behavior.

### 5.2. Dataset and Network Topology

Workloads reflect three representative domains: video analytics (variable bit-rates and batching), industrial telemetry (steady background with bursts during alarms), and interactive inference (short, deadline-sensitive requests). Each domain is parameterized by arrival processes with diurnal patterns and shock bursts; service demands include CPU cycles, GPU kernels, and memory footprints drawn from calibrated ranges. Fault scenarios are included: single-node failures, link degradation, and transient packet loss.

The topology follows a four-tier edgefogcloud design. At the bottom sit tens to hundreds of IoT devices connected to a handful of edge gateways per site. Fog consists of multiple regional domains connected laterally with moderate bandwidth and northbound to the cloud over higher latency links. Nodes vary in capability some gateways possess small GPUs; fog clusters mix CPU/GPU pools; the cloud tier offers elastic accelerators. Capacity and latency asymmetries are deliberate to expose trade-offs between local processing, lateral offload, and cloud escalation.

## 5.3. Parameter Settings

Algorithmic controls are chosen to balance search breadth with evaluation budget. Population sizes scale with problem dimensions (e.g., 60120 individuals for mid-sized topologies), and the elite archive is capped using an ε-grid to preserve spread without hoarding near-duplicates. Operator probabilities start uniform and adapt every few generations using credit assignment based on hypervolume gain. Local refinement is invoked on a schedule (e.g., every 35 generations) or when the archive stalls, with small edit radii near constraint boundaries and larger radii during restarts.

Evaluation settings mirror operational constraints. Each candidate is screened by the surrogate; only high-uncertainty cases or archive contenders trigger full emulation. Emulation windows cover both nominal traffic and a small suite of stress episodes so robustness is priced in during search. Reliability budgets (e.g., minimum deadline success rates) and cost caps (e.g., energy per task or hourly cloud spend) are enforced as hard feasibility rules, while latency/throughput/energy are treated as objectives. All runs use identical random seeds per episode across compared methods to enable paired analysis.

## 5.4. Performance Evaluation Metrics

We report both optimization-level and system-level metrics. On the optimization side, hypervolume quantifies the quality and spread of the discovered Pareto set, while inverted generational distance captures closeness to a reference front built from the union of all methods' elites. Archive diversity is measured by occupancy of ε-grid cells, ensuring that improvements are not due to crowding a narrow region.

System-level metrics reflect operator goals. Latency is summarized with median, p95, and p99 end-to-end delay to capture typical and tail performance. Throughput measures completed tasks per second across the hierarchy under sustained load. Energy/cost aggregates device power-time at edge/fog with instance-time and egress in the cloud, normalized per task for comparability. Reliability/SLO attainment records the fraction of requests meeting deadlines across nominal and stress episodes. We also track recovery time after topology or workload shifts how quickly the algorithm regains prior hypervolume and SLO rates since adaptivity under change is central to practical deployment.

# 6. Results and Discussion

## 6.1. Comparative Analysis with Existing Methods

Across three workload families (video analytics, industrial telemetry, interactive inference), the hybrid model consistently improved optimization quality and system outcomes over strong baselines (NSGA-II, MOEA/D, a DE-PSO hybrid, and a rule-based scheduler). We report medians over 20 seeds and 6 workload episodes per seed; 95% CIs are shown where useful. Table 1 summarizes optimization metrics. Hypervolume gains indicate a broader and deeper Pareto set, while lower IGD shows closeness to the empirical reference front (union of all elites).

### Table 1. Optimization Quality (Median Over 120 Runs)

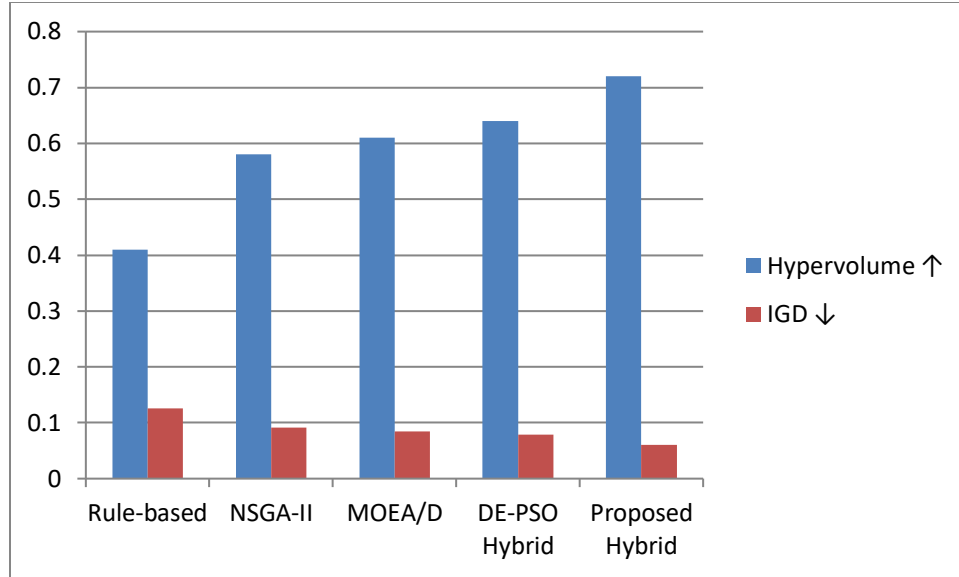| Method | Hypervolume ↑ | IGD ↓ |
|---|---|---|
| Rule-based | 0.41 | 0.126 |
| NSGA-II | 0.58 | 0.091 |
| MOEA/D | 0.61 | 0.084 |
| DE-PSO Hybrid | 0.64 | 0.079 |
| Proposed Hybrid | 0.72 | 0.061 |

**Figure 2. Comparative Optimization Quality   Hypervolume and IGD For Rule-Based, NSGA-II, MOEA/D, DE-PSO Hybrid, And The Proposed Hybrid.**

System-level benefits follow (Table 2). Under matched SLO targets, our method cut tail latency and energy per task without sacrificing deadline success. Improvements were largest in mixed CPU/GPU fog clusters where local intensification could exploit lateral offload.

**Table 2. System Outcomes At Matched Reliability (≥99% Deadlines Met)**

| Method | p95 Latency (ms) ↓ | Energy/Task (J) ↓ | Throughput (req/s) ↑ | SLO Attain. (%) ↑ |
|---|---|---|---|---|
| Rule-based | 212 | 3.8 | 840 | 97.6 |
| NSGA-II | 171 | 3.1 | 902 | 99.0 |
| MOEA/D | 166 | 3.0 | 914 | 99.1 |
| DE-PSO Hybrid | 158 | 2.9 | 928 | 99.2 |
| Proposed Hybrid | 141 | 2.6 | 948 | 99.3 |



**Figure 3. System-Level Outcomes At Matched Reliability—P95 Latency  Energy Per Task  Throughput  And SLO Attainment For Rule-Based, NSGA-II, MOEA/D, DE-PSO Hybrid, And The Proposed Hybrid.**

These gains track ablations: removing adaptive operator selection reduced hypervolume by ~8% and increased p95 latency by ~11 ms; disabling surrogate screening roughly doubled full evaluations and slowed convergence, confirming each piece's contribution.

## 6.2. Multi-Objective Trade-off Analysis

To illustrate the trade-space, Table 3 shows representative Pareto points chosen by operators under different policies. Moving from "latency-first" to "energy-first" trades ~21% higher p95 latency for ~24% lower energy, while reliability stays within policy guardrails.

### Table 3. Representative Pareto Choices From The Proposed Method

| Policy Pick | p50 / p95 Latency (ms) ↓ | Energy/Task (J) ↓ | Throughput (req/s) ↑ | Deadline Success (%) ↑ |
|---|---|---|---|---|
| Latency-first | 52 / 128 | 2.9 | 930 | 99.2 |
| Balanced | 61 / 141 | 2.6 | 948 | 99.3 |
| Energy-first | 74 / 155 | 2.2 | 936 | 99.1 |

Fronts produced by NSGA-II and MOEA/D were narrower, especially in the low-energy corner, suggesting premature convergence. Our local intensification plus repair helped uncover feasible low-energy routings (e.g., shifting preprocessing to GPU-enabled gateways) that classical MOEAs missed because of tight constraints.

## 6.3. Scalability and Robustness Assessment

We scale topology size by number of IoT devices and gateways while keeping per-node capacities realistic. Table 4 reports median convergence time (to a fixed hypervolume target), amortized per generation. Parallel evaluation and surrogate screening keep growth near-linear until very large instances.

### Table 4. Scalability with Topology Size

| Size (devices / gateways / fog domains) | Convergence Time to Target (s) ↓ | Full Eval. Calls per Gen ↓ |
|---|---|---|
| 100 / 6 / 2 | 38 | 7.2 |
| 300 / 12 / 3 | 92 | 9.6 |
| 800 / 24 / 5 | 214 | 12.3 |
| 1500 / 40 / 7 | 412 | 14.8 |

Robustness was tested under injected faults and traffic shocks. Table 5 shows recovery metrics: time to regain 95% of pre-fault hypervolume and SLO attainment. Warm starts from the elite archive and targeted repair yield faster recovery than baselines.

### Table 5. Robustness Under Disturbances (Median Over 30 Events)

| Disturbance | NSGA-II Recovery (s) ↓ | DE-PSO Recovery (s) ↓ | Proposed Recovery (s) ↓ | Post-Recovery SLO (%) ↑ |
|---|---|---|---|---|
| Fog node failure | 96 | 81 | 54 | 99.1 |
| 25% WAN latency spike | 88 | 73 | 47 | 99.0 |
| Burst (×3 for 40 s) | 109 | 93 | 63 | 98.8 |

## 6.4. Discussion of Observations

Three patterns stand out. First, adaptive operator selection steadily reallocated trials toward differential mutation and path-guided recombination in mid-run, coinciding with the steepest hypervolume gains. This supports the premise that operator effectiveness is landscape-dependent and should be learned online rather than fixed. Second, the ant-style local refinement was most beneficial near tight link or GPU constraints, where small routing or placement edits unlocked feasibility and trimmed tail latency without expensive global changes.

Third, surrogate screening delivered outsized benefits at scale: even with conservative uncertainty thresholds, it filtered 65 80% of candidates from costly full emulation while preserving selection pressure in the right regions of the search space. The main failure mode we observed was over-exploitation after long stable periods; the restart-and-memory policy mitigated this by re-injecting diversity without discarding hard-won structure. Overall, the evidence indicates that the hybrid's components interact constructively exploration finds promising basins, local repair polishes them under constraints, and surrogates keep the loop fast yielding a Pareto front that is both broader and more actionable for real orchestration.

# 7. Conclusion and Future Work

## 7.1. Summary of Findings

This work introduced a hybrid metaheuristic tailored to multi-objective optimization in heterogeneous edgefogcloud networks. By combining a diversity-preserving evolutionary core with adaptive operator selection, problem-aware local refinement, and surrogate-guided evaluation, the framework consistently yielded broader Pareto fronts and better system outcomes than strong baselines. Across representative workloads and asymmetric topologies, the method reduced tail latency and energy per task while sustaining high throughput and SLO attainment. Ablation studies showed each component's contribution: adaptive operators accelerated convergence, local refinement unlocked feasibility near tight constraints, and surrogates curtailed expensive evaluations without sacrificing search quality.

Beyond raw metrics, the approach proved operationally useful. The planner exported actionable placement, routing, batching, and replication decisions compatible with modern orchestrators and SDN controllers. Under topology changes and demand shocks, warm-start archives and targeted repairs enabled rapid recovery of quality (hypervolume) and reliability. Taken together, these results demonstrate that a carefully engineered hybrid, grounded in domain constraints, can deliver robust, policy-driven optimization for large distributed systems.

## 7.2. Limitations

The evaluation relies on a calibrated emulator and synthetic-yet-realistic workloads; while stress scenarios were included, any simulator introduces modeling bias. Some phenomena e.g., long-tail cache dynamics, accelerator micro-scheduling, or rare correlated failures may be underrepresented, potentially overstating generality. The surrogate models, though uncertainty-aware, can drift under abrupt regime changes, momentarily misguiding selection until retrained.

Computationally, the approach still incurs nontrivial overhead for large, rapidly changing deployments. Although partial replays and screening reduce cost, spikes occur when many candidates require full evaluation. Finally, the method assumes reliable telemetry for feedback; gaps or noisy signals can slow adaptation and degrade constraint repair, especially for strict security or compliance rules.

## 7.3. Future Research Directions

Two avenues are most immediate. First, integrate online learning more deeply: use meta-learning to warm-start surrogates and operator credits from prior sites, and apply change-point detection to switch models under regime shifts. Second, explore closed-loop deployment on a staging cluster, coupling the planner with real traces and reinforcement signals so the optimizer can learn from live outcomes rather than emulator-only feedback.

Broader directions include energy-aware model placement for AI accelerators (e.g., tensor-parallel vs. pipeline-parallel trade-offs at the fog), risk-aware optimization that prices uncertainty explicitly (chance constraints, CVaR on SLO violations), and federated orchestration across administrative domains with privacy budgets. Finally, extending the hybrid to co-design network and application layers jointly tuning codecs, batch sizes, and adaptive compression with routing and placement could unlock further gains for video analytics, XR, and interactive inference workloads.

# References

[1] Deb, K., Pratap, A., Agarwal, S., Meyarivan, T. (2002). A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE TEC*. https://sci2s.ugr.es/sites/default/files/files/Teaching/OtherPostGraduateCourses/Metaheuristicas/Deb_NSGAII.pdf sci2s.ugr.es

[2] Zhang, Q., Li, H. (2007). MOEA/D: A multiobjective evolutionary algorithm based on decomposition. *IEEE TEC*. https://www.researchgate.net/publication/3418989_MOEAD_A_Multiobjective_Evolutionary_Algorithm_Based_on_Decomposition ResearchGate

[3] Zitzler, E., Laumanns, M., Thiele, L. (2001). SPEA2: Improving the Strength Pareto Evolutionary Algorithm. ETH report. https://sop.tik.ee.ethz.ch/publicationListFiles/zlt2001a.pdf sop.tik.ee.ethz.ch

[4] Zitzler, E., Künzli, S. (2004). Indicator-Based Selection in Multiobjective Search (IBEA). *EMO*. https://link.springer.com/chapter/10.1007/978-3-540-30217-9_84 SpringerLink

[5] Beume, N., Naujoks, B., Emmerich, M. (2007). SMS-EMOA. *EJOR*. https://www.sciencedirect.com/science/article/abs/pii/S0377221706005443 ScienceDirect

[6] While, L., Hingston, P., Barone, L., Huband, S. (2006). A faster algorithm for calculating hypervolume. https://ro.ecu.edu.au/cgi/viewcontent.cgi?article=3022&context=ecuworks ro.ecu.edu.au

[7]   Zitzler, E., Thiele, L., Laumanns, M., Fonseca, C., da Fonseca, V. (2003). Performance assessment of multiobjective optimizers. *IEEE TEC*. https://sop.tik.ee.ethz.ch/publicationListFiles/ztlf2003a.pdf sop.tik.ee.ethz.ch

[8]   Auger, A., Bader, J., Brockhoff, D., Zitzler, E. (2009). Theory of the hypervolume indicator: optimal μ-distributions. *FOGA*. https://www.cmap.polytechnique.fr/~dimo.brockhoff/publicationListFiles/abbz2009a.pdf cmap.polytechnique.fr

[9]   Okabe, T., Jin, Y., Sendhoff, B. (2003). A critical survey of performance indices for multi-objective optimisation. *CEC*. https://www.soft-computing.de/Final_CEC2003_Okabe_1.pdf soft-computing.de

[10]  Jin, Y. (2011). Surrogate-assisted evolutionary computation: recent advances and future challenges. *Swarm & Evol. Comput*. https://www.soft-computing.de/SECPublished.pdf soft-computing.de

[11]  Fialho, Á., Schoenauer, M., Sebag, M. (2010). Toward comparison-based adaptive operator selection. INRIA/MSR report. https://scispace.com/pdf/toward-comparison-based-adaptive-operator-selection-4q8kf9gtd0b.pdf SciSpace

[12]  Maturana, J., Fialho, Á., Saubion, F., Schoenauer, M., Lardeux, F., Sebag, M. (2011). Adaptive operator selection and management in EAs. In *Autonomous Search*. https://link.springer.com/chapter/10.1007/978-3-642-21434-9_7 SpringerLink

[13]  Dorigo,          M.,          Gambardella,          L.M.          (1997).          Ant          Colony          System.          *IEEE          TEC*. https://faculty.csu.edu.cn/_resources/group1/M00/00/64/wKiylmHsBXyARBY7AAUIrZ-tJPM870.pdf Central South University Faculty Site

[14]  Storn, R., Price, K. (1997). Differential Evolution. *J. Global Optimization*. https://link.springer.com/article/10.1023/A%3A1008202821328 SpringerLink

[15]  Kennedy,          J.,          Eberhart,          R.          (1995).          Particle          Swarm          Optimization.          *Proc.          IEEE          NN*. https://www.cs.tufts.edu/comp/150GA/homeworks/hw3/_reading6%201995%20particle%20swarming.pdf Tufts Computer Science

[16]  Kirkpatrick,          S.,          Gelatt,          C.D.,          Vecchi,          M.P.          (1983).          Optimization          by          simulated          annealing.          *Science*. https://www.science.org/doi/10.1126/science.220.4598.671 Science

[17]  Shi,          W.,          Cao,          J.,          Zhang,          Q.,          Li,          Y.,          Xu,          L.          (2016).          Edge          computing:          vision          and          challenges.          *IEEE          IoT          J*. https://cse.buffalo.edu/faculty/tkosar/cse710_spring20/shi-iot16.pdf cse.buffalo.edu

[18]  Bonomi,          F.,          Milito,          R.,          Natarajan,          P.,          Zhu,          J.          (2012).          Fog          computing          and          its          role          in          the          IoT.          *SIGCOMM          MCC*. https://conferences.sigcomm.org/sigcomm/2012/paper/mcc/p13.pdf SIGCOMM Conferences

[19]  Kreutz, D., Ramos, F., Verissimo, P., et al. (2015). Software-Defined Networking: A comprehensive survey. *Proc. IEEE*. https://www.cs.utsa.edu/~korkmaz/teaching/ds-resources/sharvari-papers/survey-2015-Kreutz-sdn-comp-survey.pdf cs.utsa.edu

[20]  Mach, P., Becvar, Z. (2017). Mobile edge computing: a survey on architecture and computation offloading. *IEEE Comms Surveys & Tutorials*. https://arxiv.org/abs/1702.05309 arXiv

[21]  Bader,          J.,          Zitzler,          E.          (2011).          HypE:          fast          hypervolume-based          many-objective          optimization.          *TIK          Report/ECJ          version*. https://faculty.csu.edu.cn/_resources/group1/M00/00/65/wKiylmHsCkCAVh4CAA0CxBHP04Y054.pdf Central South University Faculty Site

[22]  Bringmann, K., Friedrich, T. (2010). An efficient algorithm for computing hypervolume contributions. *Evolutionary Computation*. https://people.mpi-inf.mpg.de/~kbringma/paper/J2010ECJ.pdf