



# INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

Volume 6, Issue 4, April 2023



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

*Impact Factor: 7.54*



6381 907 438



6381 907 438



ijmrset@gmail.com



www.ijmrset.com



# Optimizing Cloud Cost vs Reliability: Constraint-Based Scheduling for Multi-Cloud Enterprise Workloads

Madhu Babu Amarappalli

Platform Architect, Department of Cloud Architecture & Infrastructure Liberty Mutual Insurance, Mechanicsburg,  
Pennsylvania, USA

<https://orcid.org/0009-0009-2723-0256>

**ABSTRACT:** Enterprises increasingly adopt multi-cloud to reduce vendor risk and improve resiliency, yet they often face a hard tradeoff between minimizing spending and meeting strict reliability objectives. This paper proposes a **constraint-based scheduling** approach that model's workload placement as a constrained optimization problem across multiple cloud providers, regions, and service tiers. We define cost, availability, and policy constraints (e.g., data residency, latency, capacity, and affinity/anti-affinity) and solve the placement using constraint programming / integer optimization. The approach supports reliability targets expressed as availability or error-budget style objectives while explicitly accounting for cost drivers such as instance pricing and cross-cloud data movement. A practical blueprint for implementation and evaluation is presented.

**KEYWORDS:** Multi-cloud, scheduling, constraint programming, reliability, cost optimization, SLO, availability

## I. INTRODUCTION

Cloud computing enables on-demand access to configurable computing resources and has become a default platform for enterprise workloads [1]. As adoption grows, organizations increasingly operate in **multi-cloud** environments to improve resilience and reduce concentration risk, but this introduces operational complexity and cost sprawl that FinOps practices aim to control [2].

Reliability, meanwhile, is typically formalized through availability objectives (e.g., 99.9%+) and operational policies such as error budgets [3]. Major cloud architecture frameworks emphasize distributing workloads across failure domains (zones/regions) and designing for recovery and fault isolation [4], [5]. In containerized enterprise platforms (e.g., Kubernetes), schedulers already support constraints to spread replicas across failure domains for high availability [6], [7].

However, enterprise scheduling is no longer “pick the cheapest node.” It must simultaneously satisfy: (i) reliability targets, (ii) compliance and placement policies, (iii) performance/latency requirements, and (iv) budget constraints. This motivates a **constraint-based** formulation where the scheduler selects a placement that is *provably feasible* and *cost-optimal* (or *near-optimal*) under explicit reliability requirements.

## II. LITERATURE VIEW

Large-scale cluster management and scheduling have evolved from centralized systems (e.g., Borg/Omega concepts that influenced Kubernetes) to constraint-driven placement with priorities and bin-packing [8]. Kubernetes exposes practical constraint mechanisms—node selection, affinity/anti-affinity, and topology spread—to control availability and utilization [6], [7], [14].

In cloud workflow and DAG scheduling research, the cost–reliability tradeoff has been explored by multi-cloud workflow schedulers that consider reliability as a first-class objective alongside cost and makespan [9]. Surveys highlight that multi-cloud scheduling must address heterogeneous pricing, resource models, and cross-cloud data transfer constraints, and that optimization objectives often conflict (cost vs. performance vs. reliability) [10], [11].



Reliability-aware optimization has also been studied in constrained workflow settings (e.g., minimizing execution cost under makespan and reliability constraints) [12]. From an optimization tooling perspective, modern constraint programming / integer optimization solvers (e.g., CP-SAT) are designed for expressing complex discrete constraints and solving large scheduling/assignment problems efficiently [13].

### III. CONSTRAINT-BASED MULTI-CLOUD SCHEDULING MODEL AND OPERATIONALIZATION

Model and the practical elements needed to operationalize it.

#### 1) System model (workloads, options, and failure domains)

Assume an enterprise workload portfolio consisting of applications decomposed into deployable units (services, jobs, or workflow tasks). For each deployable unit  $i$ , we define a set of candidate placement options  $j$  (provider  $\times$  region/zone  $\times$  instance/service tier  $\times$  pricing model).

We treat zones/regions as failure domains and enforce replicas spreading similar to topology-aware scheduling concepts [6], [7].

#### 1.1 Workload decomposition

An enterprise workload portfolio is decomposed into *deployable units*, for example:

- **Microservices** (stateless web/API services)
- **Stateful services** (databases, caches, queues)
- **Batch jobs** (ETL, ML training, reporting)
- **Workflow tasks** in DAGs (task dependencies with data exchange)

Let the set of deployable units be  $I = \{1, 2, \dots, n\}$ . Each unit  $i$  has a specification:

- **Resource demands:**  $cpu_i, mem_i, gpu_i, disk_i$
- **Replication requirement:**  $r_i$  (e.g., 3 replicas for HA services, 1 for batch jobs)
- **SLO/reliability class:** e.g., Gold/Silver/Bronze (mapped to availability targets)
- **Compliance tags:** data residency, encryption requirements, approved providers
- **Traffic and dependency edges:**  $(i \rightarrow k)$  with expected request rate / data volume

#### 1.2 Candidate placement options

For each deployable unit  $i$ , define a set of candidate placement options:

$$J_i = \{j: (provider \times region/zone \times tier \times pricing\_model)\}$$

Each option  $j$  represents a concrete runtime target, e.g.:

- Provider = AWS / Azure / GCP
- Region = us-east / eu-west, Zone = a/b/c
- Tier = VM family, managed Kubernetes node pool, serverless tier
- Pricing model = on-demand, reserved, savings plan, spot/preemptible

For each option  $j$ , store attributes:

- Unit price  $p_j^{compute}$ ,  $p_j^{storage}$ , and transfer prices if needed
- Capacity / quota bounds  $Cap_j^{cpu}$ ,  $Cap_j^{mem}$
- Estimated reliability score  $A_{ij}$  (or  $A_j$  if option reliability is workload-agnostic)
- Compliance eligibility (allowed/blocked per policy)
- Latency/distance characteristics to other regions/providers

#### 1.3 Failure domains and spreading

Multi-cloud reliability improves when replicas are spread across **independent failure domains**:

- **Zone** (intra-region failures)
- **Region** (regional incidents)
- **Provider** (provider-wide issues, account issues, systemic outages)

Treat zones/regions as failure domains and enforce replica spreading using topology-aware scheduling concepts [6], [7]. In practice, “independence” is approximated using domain labels, for example:

- $domain(j) \in \{provider, region, zone\}$
- Topology constrains force replicas not to collapse into one domain.



**2) Cost model**

Total cost can be modeled as:

- **Computer cost:** instance/service runtime  $\times$  unit price.
  - **Storage cost:** volume/object storage  $\times$  duration.
  - **Data transfer cost:** especially cross-region/cross-cloud traffic (modeled as constraints or linear costs).
- FinOps principles motivate making these drivers explicit, measurable, and attributable to teams/services [2].

A practical cost model should match how finance actually sees bills while being “solver-friendly” (often linear or piecewise linear).

**2.1 Total cost structure**

For a planning horizon  $T$  (e.g., 1 hour, 1 day, 1 month), total cost:

$$Cost = Cost_{compute} + Cost_{storage} + Cost_{transfer} + Cost_{risk} \text{ (optional)}$$

**2.2 Compute cost**

For each deployment  $i$  assigned to option  $j$ , compute cost typically:

$$Cost_{compute} = \sum_i \sum_{j \in J_i} x_{ij} \cdot runtime_i(T) \cdot p_j^{compute}$$

Where:

- $x_{ij} \in \{0,1\}$  is the assignment decision (or per-replica decision if modeling replicas explicitly)
- $runtime_i(T)$  is expected active time during the horizon (1.0 for always-on services)

**Enterprise detail:** pricing model effects can be captured as:

- A fixed “commitment” cost for reserved capacity
- Lower marginal cost for usage under commitment
- Penalty terms for using on-demand above a threshold

If you want to keep it simple for publication, state that pricing model is embedded in  $p_j^{compute}$  per option.

**2.3 Storage cost**

Storage cost depends on volume/object size and retention:

$$Cost_{storage} = \sum_i \sum_{j \in J_i} x_{ij} \cdot storage_i(T) \cdot p_j^{storage}$$

For stateful services, including replication factor and snapshot/backup requirements.

**2.4 Data transfer cost**

Transfer is often where multi-cloud gets expensive. Model either as:

- **Linear cost:**  $volume_{ik}(T) \cdot p_{j \rightarrow l}^{transfer}$
- **Hard constraint:** forbid high-egress edges or cap cross-cloud transfer

Example:

$$Cost_{transfer} = \sum_{(i \rightarrow k)} \sum_{j \in J_i} \sum_{l \in J_k} x_{ij} x_{kl} \cdot traffic_{ik}(T) \cdot p_{j \rightarrow l}^{egress}$$

This term is quadratic. To keep it solver-friendly, operational systems often:

- Approximate using precomputed penalties per placement pair (then linearize)
- Restrict transfers with constraints (e.g., “service  $i$  and  $k$  must be co-region”)



## 2.5 FinOps operationalization

FinOps principles motivate making these drivers explicit, measurable, and attributable [2]:

- Every deployable unit has cost center tags (team/app/env)
- The scheduler produces a cost breakdown per service and per domain
- Decisions can be justified: “this placement meets SLO at lowest expected cost”

## 3) Reliability model (availability and redundancy)

Reliability targets can be expressed as:

- **Availability constraints:** e.g., service availability  $\geq$  target, aligned with standard availability interpretation in cloud architecture guidance [4].
- **Error-budget constraints:** allow a bounded level of failure/latency violation over a window, consistent with SRE practices [3].

A simple, implementable approximation is to map each placement option to an estimated availability score  $A_{ij}$  (from historical telemetry, platform SLOs, or internal reliability scoring). For replicated services, the effective availability can be approximated via redundant placement across independent failure domains (zones/regions), and enforced via constraints (e.g., at least  $k$  replicas across  $\geq 2$  zones).

Reliability requirements should be expressed as constraints that are auditable and enforceable.

## 3.1 Reliability targets as constraints

Common forms:

- **Availability constraints:**  $Availability_i \geq Target_i$  [4]
- **Error-budget constraints:** failure/latency violations bounded over a window [3]

In enterprise practice, teams often define service tiers:

- Gold: 99.95%+; multi-zone required; sometimes multi-region
- Silver: 99.9%; multi-zone preferred
- Bronze: best-effort; single-zone allowed

## 3.2 Availability scoring for placement options

A simple implementable approximation maps each option to an availability estimate:

- $A_{ij}$  = predicted availability of unit  $i$  if placed on option  $j$

This can be derived from:

- Historical incident/uptime telemetry
- Internal SLO dashboards per region/cluster
- Known risk factors (spot interruption rate, quota volatility, past outages)

You can frame  $A_{ij}$  as either:

- Direct **availability probability**, or
- A **reliability score** that ranks options consistently

## 3.3 Redundancy through replica placement

For replicated services, reliability is improved by spreading replicas across independent failure domains. In a paper, you can describe two approaches:

### Approach A (constraint-only, practical and robust):

Don't compute exact availability; enforce redundancy rules:

- At least  $k$  replicas across zones (or  $\geq 2$  regions for higher tiers)
- No two replicas in the same zone for critical services

This aligns directly with topology-aware constraints [6], [7] and multi-location guidance [4], [5].

### Approach B (approximate availability aggregation):

If you assume independence across domains, you can approximate service availability as:

$$Availability_i \approx 1 - \prod_{\text{replica } r} (1 - A_{i, \text{placement}(r)})$$



This is nonlinear; many implementations avoid exact multiplication and instead:

- Use redundancy constraints plus minimum per-replica quality thresholds
- Or linearize via piecewise approximations

### 3.4 Reliability in multi-cloud reality

Even with spreading, “independence” is imperfect (shared DNS, shared identity provider, common CI/CD). Operationally, you can add constraints like:

- Separate control-plane dependencies (e.g., avoid same IDP region)
- Require “provider diversity” for gold tier:  $\geq 2$  providers for critical services
- Limit correlated risk domains (e.g., avoid two regions with shared backbone dependency)

### 4) Constraint-based formulation (CP/MILP style)

Let binary decision variable  $x_{ij} \in \{0,1\}$  indicate whether deployable unit  $i$  is assigned to option  $j$ .

#### Core constraints (examples):

- **Assignment:** each unit placed exactly once (or exactly  $r_i$  replicas).
- **Capacity:** do not exceed quotas/capacity of each option (CPU, memory, GPUs).
- **Affinity/anti-affinity:** co-locate or separate services (e.g., separate replicas).
- **Topology spread** distribute replicas across zones/regions for HA [6], [7].
- **Compliance constraints:** data residency, restricted regions, approved providers.
- **Latency constraints:** keep dependent services within acceptable RTT bounds.
- **Reliability constraints:** enforce minimum reliability score / redundancy rules consistent with multi-location reliability guidance [4], [5].

#### Objective (single or multi-objective):

- Minimize total cost subject to reliability  $\geq$  target, **or**
- Maximize reliability subject to budget  $\leq$  limit, **or**
- Use weighted sum: minimize Cost +  $\lambda \cdot$  ReliabilityPenalty.

These problems map naturally to CP-SAT / integer optimization, which supports integer decision variables and rich constraint constructs [13].

### 4.1 Decision variables

The simplest assignment variable:

$$x_{ij} \in \{0,1\}, 1 \text{ if unit } i \text{ is assigned to option } j$$

If you model replicas explicitly:

- Let replicas be  $(i, r)$  where  $r = 1..r_i$
- Variable becomes  $x_{irj}$

### 4.2 Core constraints (examples)

#### (a) Assignment constraint

Exactly one placement per unit (or per replica):

$$\sum_{j \in J_i} x_{ij} = 1 \forall i$$

Or for replicas:

$$\sum_{j \in J_i} x_{irj} = 1 \forall i, r$$

#### (b) Capacity / quota constraint

For each option  $j$ :

$$\begin{aligned} \sum_i x_{ij} \cdot cpu_i &\leq Cap_j^{cpu} \\ \sum_i x_{ij} \cdot mem_i &\leq Cap_j^{mem} \end{aligned}$$

(Extend similarly for GPU, storage IOPS, etc.)

**(c) Affinity / anti-affinity constraints**

- Anti-affinity (separate replicas):

$$x_{irj} + x_{ir'j} \leq 1 \text{ for } r \neq r'$$

- Affinity (co-locate services  $i$  and  $k$  into same region):

$$\sum_{j \in \text{Region}(R)} x_{ij} = \sum_{j \in \text{Region}(R)} x_{kj} \quad \forall R$$

(There are several ways to express this; pick one that fits your narrative.)

**(d) Topology spread constraints (HA)**

For each service  $i$ , ensure replicas are distributed across zones/regions [6], [7]:

- “At least two zones used”
- “Max skew” type rules (no zone has too many replicas)

**(e) Compliance constraints**

If option  $j$  violates policy for service  $i$ , forbid it:

$$x_{ij} = 0 \text{ if } (i, j) \text{ is noncompliant}$$

Examples:

- Data residency: allow only  $J_i$  in specific regions
- Approved providers: remove non-approved providers from  $J_i$

**(f) Latency constraints**

For dependency edge  $i \rightarrow k$ , require region pairing within RTT threshold:

- Either forbid cross-region placement pairs that exceed latency
- Or cap the number of “far edges” for performance-critical paths

**(g) Reliability constraints**

Several implementable patterns:

- Minimum placement quality:

$$\sum_{j \in J_i} x_{ij} \cdot A_{ij} \geq A_i^{\min}$$

- Redundancy rule constraints (recommended in practice):
  - $\geq 2$  zones
  - $\geq k$  replicas
  - Optional provider diversity for top-tier services

This aligns with multi-location reliability guidance [4], [5].

**4.3 Objective functions**

Common objectives:

**(1) Minimize cost subject to reliability constraints**

$$\min \text{ Costs.t. constraints and } Reliability_i \geq Target_i$$

**(2) Maximize reliability subject to a budget**

$$\max ReliabilityScores.t. Cost \leq Budget$$

**(3) Weighted tradeoff**

$$\min(Cost + \lambda \cdot ReliabilityPenalty)$$

Where penalty might represent:

- Violations of soft constraints
- Expected downtime cost
- Risk score for correlated domains

These problems map naturally to CP-SAT / integer optimization frameworks supporting rich constraints [13].



### 5) Practical scheduling workflow (enterprise-ready)

A pragmatic implementation loop:

1. **Ingest requirements:** per-service SLO/availability class, compliance tags, resource requests.
2. **Generate options:** candidate clouds/regions/tiers; filter forbidden placements.
3. **Score inputs:** cost estimates, risk/reliability scores, latency estimates.
4. **Solve:** run CP-SAT to produce a feasible, optimized plan [13].
5. **Enforce in runtime:** translate decisions into deployment policies (e.g., Kubernetes constraints) [6], [7].
6. **Observe + reoptimize** periodically re-solve when prices, demand, or reliability posture changes.

A production-grade system needs more than a solver—it needs a reliable pipeline of inputs, governance, and continuous re-optimization.

#### 5.1 Ingest requirements (policy + SLO + resources)

Inputs typically come from:

- Service catalog (ownership, tier, dependencies)
- IaC manifests (CPU/memory requests, replica counts)
- SRE policy (SLO targets, error budgets) [3]
- Security/compliance policy engine (region/provider restrictions)

Output: a normalized specification per deployable unit.

#### 5.2 Generate options (candidate clouds/regions/tiers)

Construct  $J_i$  by:

- Enumerating allowed providers/regions
- Enumerating feasible compute tiers (meets CPU/mem/GPU)
- Applying quota filters and “blocked” lists
- (Optional) adding “preferred” pools (existing reserved capacity)

This stage often removes 80–95% of theoretical options.

#### 5.3 score inputs (cost, reliability, latency)

Maintain three continuously updated datasets:

- **Pricing tables** (compute/storage/transfer)
- **Reliability/incident history** per provider-region-zone/cluster
- **Latency matrix** (region-to-region RTT estimates, or measured service mesh telemetry)

Convert them into solver-ready scalars:

- $p_j$  cost coefficients
- $A_{ij}$  availability/reliability scores
- Allowed/forbidden pairing tables for latency/compliance

#### 5.4 Solve (CP-SAT)

Run CP-SAT to produce:

- A feasible placement plan
- Objective value and constraint satisfaction report
- (Optional) alternatives: top-k solutions for human review [13]

Operational detail: large portfolios are usually solved by:

- Partitioning by environment (prod vs dev), geography, or dependency clusters
- Using rolling horizon solving (solve what changes, keep stable assignments fixed)

#### 5.5 Enforce at runtime (Kubernetes / platform)

Translate solver decisions into enforceable policies:

- Node labels/taints and **node selectors**
- Affinity/anti-affinity rules
- Topology spread constraints for HA [6], [7]
- Admission control / policy-as-code to prevent drift

This ensures the “plan” becomes real scheduling behavior.





### 5.6 Observe + reoptimize (closed-loop control)

Preoptimization triggers:

- Price changes (reserved/spot shifts, increases)
- Reliability posture changes (incident spikes in a region)
- Demand growth (capacity pressure)
- Policy updates (new compliance requirements)

A typical cadence:

- **Daily** for cost optimization in non-critical environments
- **Weekly** or event-driven for production to avoid unnecessary churn

To limit instability, introduce **migration cost** or **change budget** constraints (e.g., “move at most 5% of services per run unless incident-driven”).

## IV. CHALLENGES

1. **Reliable reliability inputs:** Availability/risk scoring can be noisy and non-stationary (incidents, provider outages, changing dependencies). Using simple scores is practical but imperfect; better models may be needed for mission-critical workloads [3], [4].
2. **Cross-cloud data gravity and egress:** Data transfer costs and latency can dominate, and constraints can become non-linear.
3. **Heterogeneity:** Providers differ in instance families, networking, and managed service semantics; mapping them to a unified option set is complex [11].
4. **Scale and solver performance:** Real portfolios may include thousands of services and constraints; decomposition, heuristics, and incremental solving may be required even with strong solvers [13].
5. **Operational alignment:** Scheduling decisions must align with SRE/FinOps governance—error budgets, change management, and accountability models [2], [3].
6. **Constraint drift:** Policies evolve (security, compliance, org rules). Ensuring the model matches reality is an ongoing engineering task.

## V. FUTURE

- **Stochastic / robust optimization:** explicitly model uncertainty in demand, outages, and price fluctuations to avoid brittle placements.
- **Closed-loop reliability control:** integrate error-budget burn rate signals to trigger re-placement or replica adjustments [3].
- **Multi-objective Pareto planning:** present decision-makers with Pareto frontiers (cost vs. reliability vs. latency) rather than a single answer.
- **Deeper platform integration:** generate Kubernetes-native policies (affinity/spread/bin-pack tuning) directly from the solver output [6], [14].
- **Governance automation:** tie outputs to FinOps tagging, show back/chargeback, and policy-as-code enforcement [2].

## VI. CONCLUSION

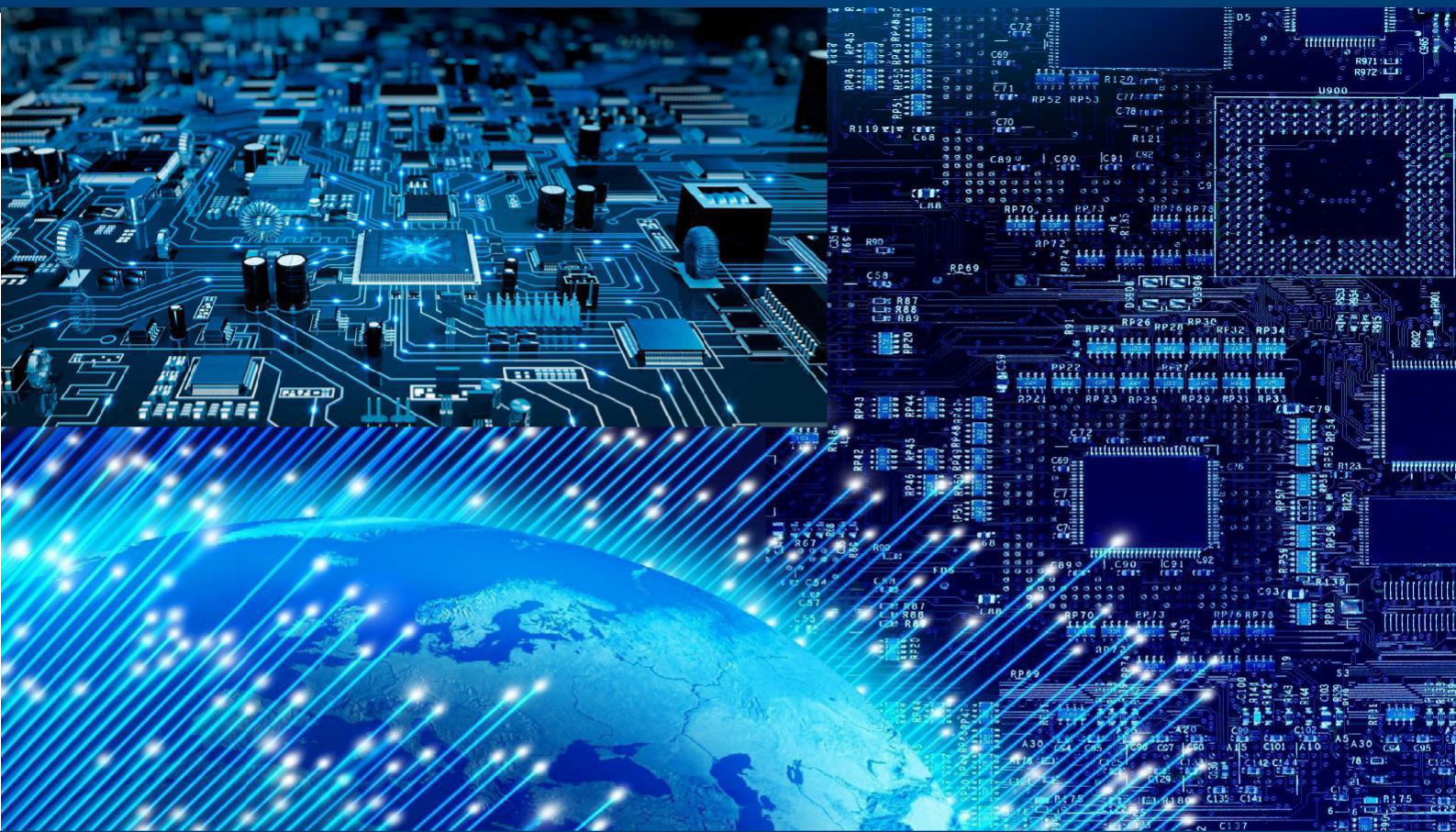
Multi-cloud enterprise scheduling requires balancing cost control with reliability guarantees under complex real-world constraints. This paper presented a constraint-based scheduling approach that formalizes placement as an optimization problem with explicit reliability objectives, compliance rules, and topology constraints. By leveraging modern constraint programming / integer optimization and enforcing results through platform schedulers, enterprises can produce placements that are both economically efficient and reliability-aware, while remaining auditable and policy-driven.

## REFERENCES

1. P. Mell and T. Grance, “The NIST Definition of Cloud Computing (SP 800-145),” NIST, 2011. [NIST Publications](#)
2. FinOps Foundation, “FinOps Framework / What is FinOps,” updated Jan. 2021. [FinOps Foundation+1](#)
3. Google SRE, “Error Budget Policy,” *Site Reliability Engineering Workbook*. [Google SRE](#)



4. AWS, “Well-Architected Reliability Pillar: Deploy the workload to multiple locations (REL10-BP01),” AWS Documentation. [AWS Documentation](#)
5. Microsoft, “Azure Well-Architected Framework — Reliability,” Microsoft Learn. [Microsoft Learn+1](#)
6. Kubernetes, “Pod Topology Spread Constraints,” Kubernetes Documentation (updated Oct. 27, 2022). [Kubernetes](#)
7. Kubernetes, “Assigning Pods to Nodes,” Kubernetes Documentation (updated Aug. 2, 2021). [Kubernetes](#)
8. B. Burns et al., “Borg, Omega, and Kubernetes,” (paper PDF). [Google Research](#)
9. A. Tekawade and S. Banerjee, “Cost and Reliability Aware Scheduling of Workflows Across Multiple Clouds with Security Constraints,” arXiv:2304.00313, 2021. [arXiv](#)
10. M. U. Sana et al., “Efficiency aware scheduling techniques in cloud computing,” 2021 (open-access article). [PubMed Central](#)
11. Q. Zhang et al., “Survey on Task Scheduling Optimization Strategy under Multi-Cloud,” 2022. [ScienceDirect](#)
12. X. Fu et al., “Reliability Aware Cost Optimization for Cloud Workflows under Constraints,” (paper PDF). [East China Normal University Faculty](#)
13. Google, “CP-SAT Solver — OR-Tools,” Google Developers Documentation (updated Aug. 28, 2022). [Google for Developers](#)
14. Kubernetes, “Resource Bin Packing,” Kubernetes Documentation (updated Oct. 10, 2023). [Kubernetes](#)
15. Naga Ramesh Palakurti, [Governance Strategies for Ensuring Consistency and Compliance in Business Rules Management](#), 2023, <https://philpapers.org/rec/NAGGSF>



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | [ijmrset@gmail.com](mailto:ijmrset@gmail.com) |

[www.ijmrset.com](http://www.ijmrset.com)