

Why Cloud Cost Optimization is Important

1. Cloud Wastage is Real

Most companies **overpay** for cloud services — often by 30% to 40% — without realizing it. Why?

- Idle servers running overnight
- Overprovisioned instances
- Orphaned storage or IPs
- Expensive instance types used for basic workloads

The Pizza Shop Analogy

- Imagine you're running a pizza shop:
- You hire 10 chefs even though only 2 are cooking at a time
- You keep the **oven running 24/7**, even at night
- You store expired ingredients that you'll never use
- You buy the most expensive cheese, even for simple breadsticks
- This is exactly how companies treat cloud resources if they don't optimize costs:
- Unused servers = idle chefs
- Overprovisioned CPUs = expensive cheese
- No auto-shutdown = oven running at midnight
- Unused disks = expired ingredients
- Would you call that a smart business decision? Probably not.
 Same applies to cloud spending.



How It Helps the Company / Team

- ✓ 1. Cost Savings = Higher Profit
- Every ₹1 saved on cloud is a **direct saving** for the business
- It frees up budget for better tools, performance improvements, or hiring
- **2.** Scalable and Sustainable Cloud Usage
- Teams can scale confidently knowing they're not overpaying
- Enables better architecture planning for growth
- **3. Team Accountability**
- Helps allocate costs per team/project (e.g., via tagging)
- Teams become more responsible and aware of their infrastructure decisions

- How It Helps You as a DevOps Engineer
- 1. Career Growth & Recognition
- Teams notice when you save lakhs in cloud bills
- You become the "go-to" person for efficiency and reliability
- Gives you visibility to engineering managers, architects, and finance teams
- 2. Learn FinOps a High-Demand Skill
- Cloud + Finance = **FinOps** one of the most in-demand DevOps cross-skills
- Understanding costs shows that you're not just a tool-user but a business-aware engineer
- **3. Better Architecture Thinking**
- Cost optimization pushes you to design leaner, smarter, scalable systems
- You'll naturally improve in cloud-native design, serverless, IaC, automation, and monitoring

6 1. Rightsizing Compute Resources

What It Means:

Choosing the **right instance size** (vCPU, memory, etc.) for your workload. Many engineers overprovision VMs to "be safe" — which wastes money.

Examples:

- •You run a t3.2xlarge instance for a dev server that's idle 90% of the time. It can be downgraded to t3.medium.
- •Kubernetes pods request 2 CPU but never go above 0.3 CPU scale down requests and limits.

Tools:

•AWS: Compute Optimizer, Trusted Advisor

•Azure: Advisor

•GCP: Recommender

•Kubernetes: Goldilocks, Kube-resource-report

Tips:

- Review CPU/memory usage with monitoring tools (Prometheus + Grafana, CloudWatch).
- Use Auto Scaling Groups with dynamic instance types (t3, t4g, m6i, etc.)
- Consider ARM-based processors (AWS Graviton2/3).

Sheduling auto-Shutdown of Non-Prod + resources Scheduling and Auto-ducoioweof hen resources to save cost

© 2. Scheduling & Auto-Shutdown of Non-Prod Resources

What It Means:

 Automatically stop idle or unused resources especially dev/test/staging during off-hours.

Examples:

- Shut down staging servers at 8 PM and start them at 8 AM.
- Stop Kubernetes clusters on weekends if they aren't needed.

🔪 🦴 Tools:

- **AWS**: Instance Scheduler, Lambda + EventBridge
- Azure: Automation Runbooks
- GCP: Cloud Scheduler + Cloud Functions
- K8s: CronJobs + scripts to scale to 0

Tips:

- Use tagging like Environment=Dev and schedule shutdowns for those.
- Set up reminders or slack alerts before shutdown.

6 3. Use Spot / Preemptible / Low-Priority VMs

What It Means:

Use spare capacity from cloud providers at **up to 90% discount**, with the risk of interruptions.

Examples:

- •Use **EC2 Spot Instances** for Jenkins agents or batch jobs.
- •Run **Spark** or ML jobs on Spot/Preemptible instances.

♦ Tools:

- •AWS: EC2 Spot Fleet, Auto Scaling with Mixed Instances
- •GCP: Preemptible VMs, MIGs (Managed Instance Groups)
- •K8s: Use taints/tolerations or node selectors to isolate workloads on Spot nodes



- · 🔽 Tips:
- Always have fallbacks (On-Demand fallback policy).
- Avoid using Spot for stateful or critical apps.

6 4. Storage Optimization

What It Means:

Avoid paying for unused or over-retained data.

Examples:

- •Delete unattached EBS volumes and unused snapshots.
- •Move logs older than 30 days to **Glacier Deep Archive** or **Azure Cool Storage**.

Tools:

•AWS: S3 Lifecycle Policies, EBS Snapshot Manager, Trusted

Advisor

•Azure: Blob lifecycle rules

•GCP: Object lifecycle management



- **Tips:**
- Automate cleanup using Boto3, GCP SDK, Azure CLI.
- Visualize storage usage with CloudWatch dashboards or Grafana.

6 5. Container Resource Optimization

What It Means:

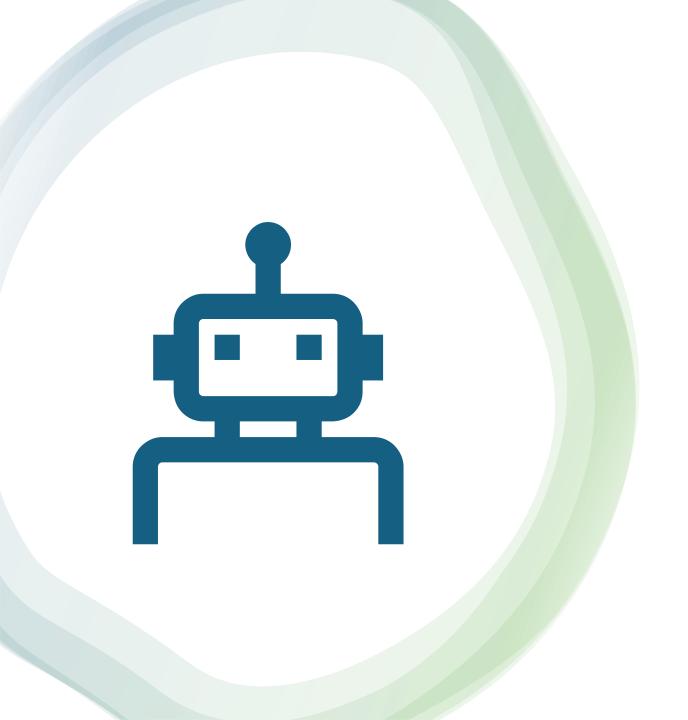
Optimize container resource requests/limits and reduce Kubernetes node waste.

Examples:

- •If a pod requests 2 vCPU but uses only 0.3 consistently, the remaining 1.7 CPU is **wasted**.
- •Orphaned PVCs or idle services still consume storage/networking resources.

↑ Tools:

- •KubeCost, Goldilocks, VPA (Vertical Pod Autoscaler)
- Prometheus + custom dashboards



- **Tips:**
- Use **Vertical Pod Autoscaler** in recommendation mode.
- Delete unused namespaces, ConfigMaps, PVCs.

6. Serverless Where It Makes Sense

What It Means:

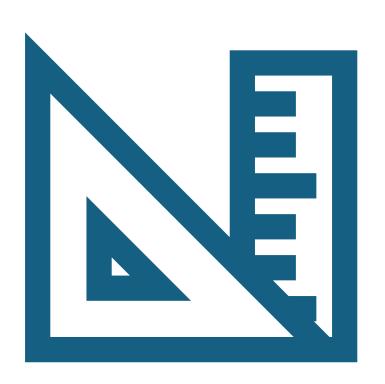
Shift from always-on servers to event-driven, pay-perexecution compute.

Examples:

- •Replace EC2 cron job with AWS Lambda (₹0 cost if idle).
- •Use Azure Functions for email processing, Cloud Functions for simple APIs.

Platforms:

- •AWS Lambda, GCP Cloud Functions, Azure Functions
- •AWS Fargate (serverless containers)



☑ Tips:

- Ensure cold-start isn't a problem for your workload.
- Use **step functions** for orchestration.

6 7. Region & AZ-Aware Cost Planning

- What It Means:
- Some regions and AZs are cheaper than others.
- § Examples:
- Running EC2 in us-east-1 is ~15-20% cheaper than ap-south-1.
- GCP Tokyo vs. Singapore cost difference for same machine type.
- \ \ Tools:
- Cloud pricing calculators (AWS, GCP, Azure)
- Terraform locals to parameterize region

- **Tips:**
- For latency-insensitive workloads (backup, batch), choose cheapest region.
- Consider **cross-region** replication cost.

6 8. Monitoring & Budget Alerts

What It Means:

Track spend proactively, not reactively.

Examples:

- •Set monthly budgets for projects and alert via Slack/email if 80% used.
- •Detect anomalies like cost spikes in a region.

Tools:

•AWS: Budgets, Cost Explorer, Cost Anomaly Detector

•Azure: Cost Management + Budget Alerts

•GCP: Billing Alerts, BigQuery for billing exports



- Use tagging to separate costs per team/project.
- Create daily burn rate dashboards